

Week 1 Tues

quantify faithfulness / closeness

average each individual data point
 ↴ ↴ ↴

Define loss function

$$\text{loss}(\bar{x}; x_i) = (\bar{x} - x_i)^2$$

← squared loss/error

Average loss function

$$R_{\text{sq}}(\bar{x}; D) = \frac{1}{n} \sum_{i=1}^n \text{loss}(\bar{x}; x_i) = \frac{1}{n} \sum_{i=1}^n (\bar{x} - x_i)^2$$

↑
 empirical risk

Find an estimator h to optimize empirical risk?

$\arg \min$: argument h that minimizes
 f(h^*) $\leq f(h)$, $\forall h \in \mathbb{R}$

$$\bar{x} = \arg \min_{h \in \mathbb{R}} R_{\text{sq}}(h; D).$$

$$= \arg \min_{h \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n (h - x_i)^2$$

empirical risk minimization

$$h^* = \frac{1}{n} \sum_{i=1}^n x_i = \text{Mean}(D)$$

The best prediction, in terms of mean squared error, is the **mean** (unique)

more points example

$$2 \text{ point example: } \frac{\partial}{\partial h} \frac{1}{2} [(h - x_1)^2 + (h - x_2)^2]$$

$$= \frac{1}{2} \frac{\partial}{\partial h} ((h - x_1)^2 + (h - x_2)^2)$$

$$= \frac{1}{2} \left(\frac{\partial}{\partial h} (h - x_1)^2 + \frac{\partial}{\partial h} (h - x_2)^2 \right)$$

$$= \frac{1}{2} (2(h - x_1) + 2(h - x_2))$$

$$= \frac{1}{2} (2(h - x_1) + h - x_2))$$

$$= 2h - x_1 - x_2$$

||

$$2h - x_1 - x_2 = 0$$

$$h = \frac{x_1 + x_2}{2}$$

$$\frac{\partial}{\partial h} \left(\frac{1}{n} \sum_{i=1}^n (h - x_i)^2 \right)$$

$$= \frac{1}{n} \left(\frac{\partial}{\partial h} \sum_{i=1}^n (h - x_i)^2 \right)$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial h} (h - x_i)^2$$

$$= \frac{1}{n} \sum_{i=1}^n 2(h - x_i)$$

$$= \frac{2}{n} \sum_{i=1}^n (h - x_i)$$

$$= \frac{2}{n} (nh - \sum_{i=1}^n x_i)$$

$$\frac{2}{n} (nh - \sum_{i=1}^n x_i) = 0$$

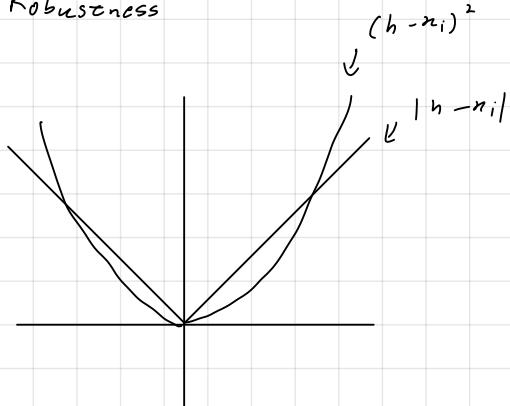
$$nh - \sum_{i=1}^n x_i = 0$$

$$\bar{x} = h^* = \frac{1}{n} \sum_{i=1}^n x_i$$

Outlier

Mean is sensitive to outliers.

Robustness



$$R^{ab}(h; D) = \frac{1}{n} \sum_{i=1}^n |h - x_i|$$

mean error
(good for robustness)

Week 1 Thurs

For dataset $D = \{x_1, \dots, x_n\}$

$\text{loss}(h, x_i)$

Empirical risk : $R(h, D) = \frac{1}{n} \sum_{i=1}^n \text{loss}(h, x_i)$.

Optimize: $\underset{h \in R}{\operatorname{argmin}} R(h, D) = h^*$

Test : $R(h^*, D_{\text{test}})$.

Def. A loss function $L(h, x)$ takes in a prediction h and a true value, x , and outputs a number measuring how far h is from x

loss function:

- Squared loss - mean (D) - $\text{Loss}_{sq}(h, x_i) = (h - x_i)^2$
- Sensitive to outliers.

- Absolute loss : $\text{loss}_{abs}(h, x_i) = |h - x_i|$
- Median (D)

optimize: $\underset{h \in R}{\operatorname{argmin}} R_{abs}(h, D) = \frac{1}{n} \sum_{i=1}^n \text{loss}_{abs}(h, x_i)$

$$\lim_{\Delta \rightarrow 0} R_{abs}(h + \Delta, D) = R_{abs}(h, D)$$

$$\frac{\partial}{\partial h} R_{abs}(h, D) = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial h} \text{loss}(h, x_i).$$

$$\mathbb{1}\{h < x_i\} = \begin{cases} 1 & , h < x_i \\ 0 & , h \geq x_i \end{cases}$$



$$\frac{\partial}{\partial h} \text{loss}(h, x_i) = \frac{\partial}{\partial h} |h - x_i| = \begin{cases} 1 & h > x_i \\ 0 & h = x_i \\ -1 & h < x_i \end{cases}$$

$$\approx -1 \cdot \mathbb{I}\{h < x_i\} + 1 \cdot \mathbb{I}\{h > x_i\}$$

$$\Rightarrow \frac{\partial}{\partial h} R_{\text{abs}}(h, D) = \frac{1}{n} \sum_{i=1}^n (-1 \cdot \mathbb{I}\{h < x_i\} + 1 \cdot \mathbb{I}\{h > x_i\})$$

$$= \frac{1}{n} \left[-1 \sum_{i=1}^n \mathbb{I}\{h < x_i\} + (1) \sum_{i=1}^n \mathbb{I}\{h > x_i\} \right]$$

$\sum_{i=1}^n \mathbb{I}\{h < x_i\}$ = # of x greater than h .

$\sum_{i=1}^n \mathbb{I}\{h > x_i\}$ = # of x less than h .

$h^* = \text{Median}(D)$,

best prediction for mean absolute error.

$$= \frac{1}{n} \left[\underbrace{\# \text{ of } x \text{ less than } h}_{x_i < h} - \underbrace{\# \text{ of } x \text{ greater than } h}_{x_i > h} \right]$$

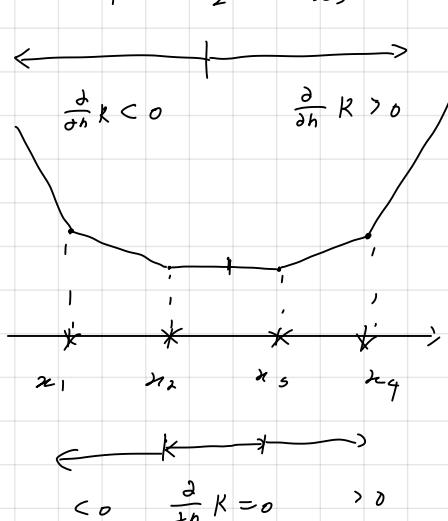
odd:
number is unique

even:
any number between
the middle two data
point minimizes mean
absolute error.

loss function (cont.)

$$\text{loss}(h, x_i) = \mathbb{I}\{h \neq x_i\}$$

ERM \rightarrow Mode(D)



$$\text{Ex. } R_{\text{abs}}(h, D) = \frac{1}{n} \sum_{i=1}^n |h - x_i|$$

$$R_{\text{sq}}(h, D) = \frac{1}{n} \sum_{i=1}^n (h - x_i)^2$$

$$\sqrt{R_{\text{sq}}(h, D)}$$

$$\text{If } R_{\text{abs}}(h, D) \leq \varepsilon, \Rightarrow \sqrt{R_{\text{sq}}(h, D)} \leq \varepsilon$$

Empirical risk — the average loss on the data sets:

$$R_L(h) = \frac{1}{n} \sum_{i=1}^n L(h, y_i)$$

T

R = "risk"

The goal of learning: find h that minimizes R_L .

This is called empirical risk minimization (ERM).

- The choice of loss function determines the properties of the result.

↓
different loss function
" " minimizer
" " prediction

ex). $R_{0,1}(h) = \frac{1}{n} \sum_{i=1}^n \begin{cases} 0 & \text{if } h = y \\ 1 & \text{if } h \neq y \end{cases} \Rightarrow h^* = \text{Mode}(P)$

loss	Minimizer	outliers	Differentiable
Labs	median	insensitive	no
Lsq	mean	sensitive	yes
$R_{0,1}$	mode	insensitive	no

$$\text{Labs}(h, y) = |y - h|$$

$$\text{Lsq}(h, y) = (y - h)^2$$

$$\text{Empirical risk } R(h) = \frac{1}{n} \sum_{i=1}^n \text{Loss}(h, y).$$

$$h^* = \arg \min_{h \in \mathcal{H}} R(h; \theta).$$

$$\begin{aligned} \bullet R_{\text{abs}}(h^*; b) &= \frac{1}{n} \sum_{i=1}^n |h^* - x_i| \\ &\approx \frac{1}{n} \sum_{i=1}^n |\text{Median}(b) - x_i| \end{aligned}$$

$$R(h^*; b)$$

$$\begin{aligned} R_{\text{sq}}(h^*; b) &= \frac{1}{n} \sum_{i=1}^n (h^* - x_i)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (\bar{x} - x_i)^2 \quad \text{Variance} \end{aligned}$$

$$D = \{1, 2, 2, +\}$$

$$R_{\text{ab}}(h^*; b) = \frac{1}{2}$$

$$\begin{aligned} \bullet R_{0,1}(h; b) &= \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{h \neq x_i\} \\ &= \frac{n-1}{n} \end{aligned}$$

$$\text{Standard deviation } \sqrt{R_{\text{sq}}(h^*; b)}$$

$$\underbrace{K_{sq}(h; b)}_{\Downarrow} \quad \text{vs.} \quad K_{abs}(h; b)$$

$$\begin{aligned}
 K_{sq}(h; b) &= \frac{1}{n} \sum_{i=1}^n (h - x_i)^2 \\
 &= \frac{1}{n} \sum_{i=1}^n |h - x_i|^2 \\
 &\stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n |h - x_i| = a_i \\
 &= \frac{1}{n} \sum_{i=1}^n a_i^2 \\
 &\stackrel{\text{def}}{=} \frac{1}{n} \sum_{j=1}^n \left(\frac{1}{n} \sum_{i=1}^n a_i^2 \right) \\
 &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n a_i^2
 \end{aligned}$$

$$\begin{aligned}
 K_{sq}(h; b) - K_{abs}(h; b) &= \frac{1}{n^2} \sum_i \sum_j \left(\frac{a_i^2 + a_j^2}{2} - a_i a_j \right) \\
 &= \frac{1}{n^2} \sum_i \sum_j \left(\frac{1}{2} (a_i^2 + a_j^2 - 2a_i a_j) \right) \geq 0 \\
 K_{sq}(h; b) &\geq K_{abs}(h; b)
 \end{aligned}$$

$$\sqrt{K_{sq}(h; b)} \geq K_{ab}(h; b), \forall h.$$

Note: $\frac{a^2 + b^2}{2} \geq ab$

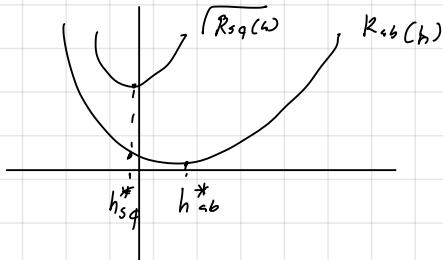
Young's Ineq. S. (a)

want:

$$\sqrt{K_{sq}(h^*; b)} \quad \text{vs.} \quad K_{ab}(h^*; b).$$

$$\min \sqrt{K_{sq}(h; b)} \geq \min K_{ab}(h; b)$$

$$\begin{aligned}
 K_{abs}(h_{ab}^*, b) &\leq K_{ab}(h_{sq}^*, b) \\
 &\leq \sqrt{K_{sq}(h_{sq}^*, b)}
 \end{aligned}
 \Rightarrow$$



Center and Spread

- The input h^* that minimizes $R(h)$ is some measure of the center of the data set.
e.g. median, mean, mode.
- The minimum output $R(h^*)$ represent some measure of the spread, or variance, in the data set.

↓

For absolute loss: $R_{\text{abs}}(h^*) = R_{\text{abs}}(\text{Median}(y_1, y_2, \dots, y_n))$

$$\uparrow \quad = \frac{1}{n} \sum_{i=1}^n |y_i - \text{median}(y_1, y_2, \dots, y_n)|.$$

this is the mean absolute deviation from the median.

For squared loss: $R_{\text{sq}}(h^*) = R_{\text{sq}}(\text{Mean}(y_1, \dots, y_n))$

$$= \frac{1}{n} \sum_{i=1}^n (y_i - \text{Mean}(y_1, \dots, y_n))^2$$

↑
this is called the variance

↑
the square root of variance is standard deviation.

Huber loss / Gradient Descent

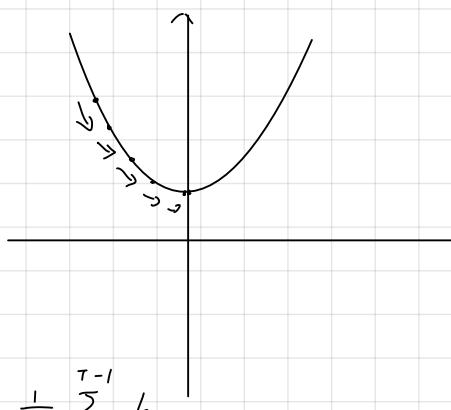
$$\text{loss}(h, x_i) = \begin{cases} |h - x_i| - \frac{1}{2}, & |h - x_i| > \frac{1}{2} \\ (h - x_i)^2, & |h - x_i| \leq \frac{1}{2} \end{cases}$$

$$R_{\text{Huber}} = \frac{1}{n} \sum_{i=1}^n \text{loss}(h, x_i)$$

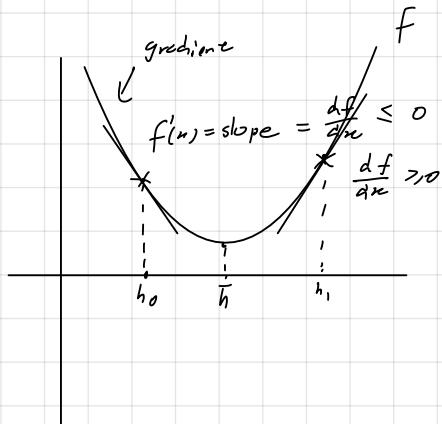
Initialize $h_0 = 0$

$$h_{t+1} = h_t - \Delta_t \cdot \frac{d}{dh} R_{\text{Huber}}(h, 0)$$

$$\frac{d h(t)}{d t} = \frac{\partial}{\partial h} R_{\text{Huber}}(h_t, 0) \quad \text{Output } h_t = \frac{1}{t} \sum_{\tau=1}^{t-1} h_\tau$$



Gradient Descent



From Math 20c:

gradient \rightarrow direction increase the most

- Pick α to be a positive number, the learning rate.

- Pick a starting prediction, h_0 .

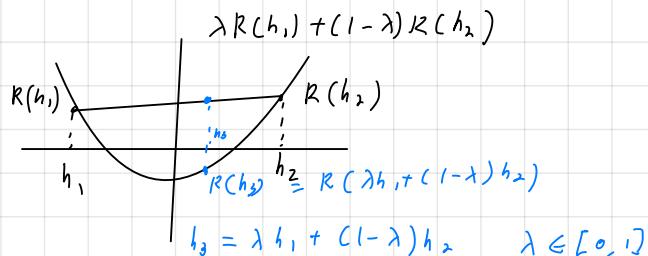
- On step i , perform update $h_i = h_{i-1} - \alpha \cdot \frac{dR}{dh}(h_{i-1})$

- Repeat until convergence.

$$x_{t+1} = x_t - \Delta \frac{df(x_t)}{dx}$$

$t = 0, 1, 2, \dots$

Convex function



$$R(h_3) = R(\lambda h_1 + (1-\lambda) h_2)$$

$$R(\lambda h_1 + (1-\lambda) h_2) \leq \lambda R(h_1) + (1-\lambda) R(h_2)$$

$$\forall \lambda \in [0, 1], \forall h_1, h_2,$$

(5)b. loss convex $\Rightarrow R(h, \theta)$ convex.

Thm: For Huber loss, with suitable $\{\Delta_0, \dots, \Delta_T\}$

step size

\nwarrow

this means hubber loss will converge

$\bar{h} = \frac{1}{T} \sum_{t=0}^{T-1} h_t$

$R(\bar{h}_T, D) - R(h^*, D) \leq \frac{4(h_0 - h^*)^2}{T}$

all h_t from gradient descent algorithm

number of iterations that G.D. converges.

If you want error ε :

$$\text{Need } T \geq \frac{4(h_0 - h^*)^2}{\varepsilon} = \frac{4(h^*)^2}{\varepsilon} \leq \frac{4K^2}{\varepsilon}$$

$$\Delta t = \frac{1}{4}, \quad \forall t. \quad \text{for Huber}$$

range of data point

$$\text{proof: } \frac{\partial^2}{\partial h^2} \text{loss}(h, h_i) = \begin{cases} 0, & |h - h_i| > \frac{1}{2} \\ 2, & |h - h_i| \leq \frac{1}{2} \end{cases} \geq 0$$

Huber loss
only has minimum,
no maximum

$$0 \leq \frac{\partial^2}{\partial h^2} \text{loss}(h, h_i) \leq 2$$

↑
convex smooth function
characteristic

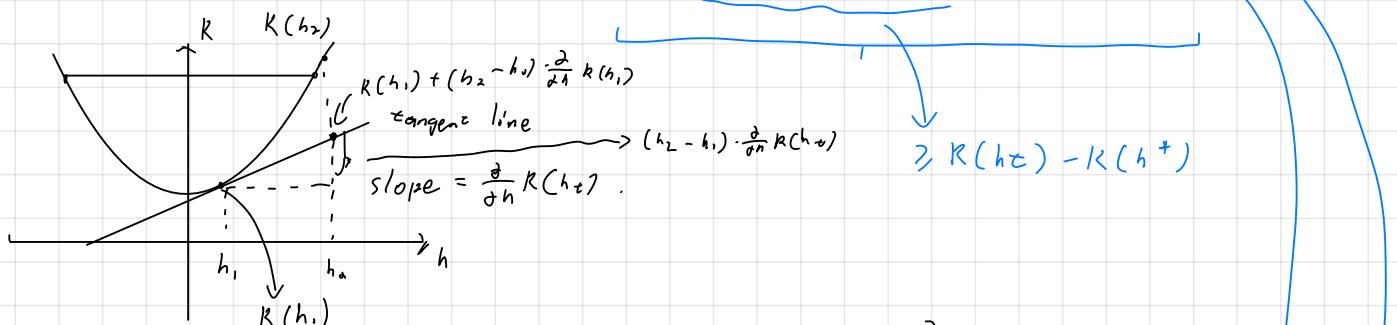
$$\Rightarrow 0 \leq \frac{\Delta^2 K(h, b)}{\Delta h^2} \leq 2$$

$$\Rightarrow \left(\frac{\partial}{\partial h} K(h) \right)^2 \leq 4(K(h) - K(h^*)) \quad \text{for } h \neq h^*.$$

$$(h_{t+1} - h^*)^2 = (h_t - \Delta t \cdot \frac{\partial}{\partial h} K(h_t) - h^*)^2$$

$$= [(h_t - h^*) - \Delta t \cdot \frac{\partial}{\partial h} K(h_t)]^2$$

$$= (h_t - h^*)^2 - 2\Delta t (h_t - h^*) \cdot \frac{\partial}{\partial h} K(h_t) + \Delta t^2 \left(\frac{\partial}{\partial h} K(h_t) \right)^2$$



$$\text{So, } K(h_2) \geq K(h_1) + (h_2 - h_1) \cdot \frac{\partial}{\partial h} K(h_1)$$

$$\Rightarrow (h_2 - h_1) \cdot \frac{\partial}{\partial h} K(h_1) \geq K(h_1) - K(h_2)$$

$$(h_t - h^*) \cdot \frac{\partial}{\partial h} K(h^*) \geq K(h_t) - K(h^*)$$

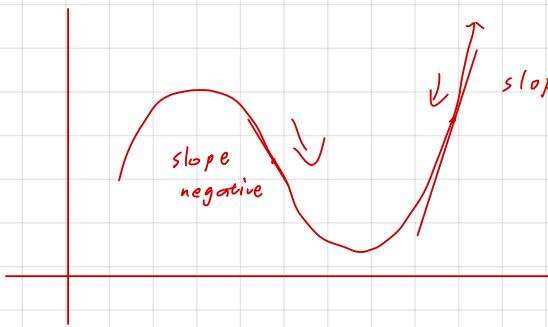
$$\leq (h_t - h^*)^2 - 2\Delta t (K(h_t) - K(h^*)) + \Delta t^2 (K(h_t) - K(h^*))$$

$$\leq (h_t - h^*)^2 - 2\Delta t (1 - 2\Delta t) (R(h_t) - R(h^*)).$$

maximize

$$\Psi \\ \Delta t = \frac{1}{4}$$

Gradient Descent



slope is positive initial guess.

$$h_1 = h_0 - \frac{dR}{dh} / \nabla R(h_0)$$

opposite the direction of
the derivative / gradient.

Algorithm :

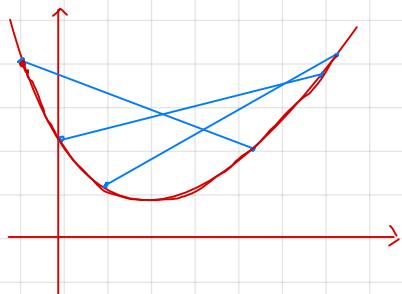
- Pick α to be positive number, the learning rate, also known as step size.
- Pick a starting prediction, h_0 .
- On step i , perform update $h_i = h_{i-1} - \alpha \cdot \frac{dR}{dh}(h_{i-1})$
- Repeat until convergence.

Convex function

f is convex if for every a, b in the domain of f , the line segment between

$$(a, f(a)) \text{ & } (b, f(b)).$$

does not go below the plot of f .



function $f : K \rightarrow K$ is convex if $\forall a, b \in K$ and $t \in [0, 1]$,

$$(1-t)f(a) + t f(b) \geq f((1-t)a + tb).$$

Thm: If $R(h)$ is convex and differentiable, then gradient descent converges to a global minimum of R provided that the step size is small enough.

↑

because f convex and has local min \Rightarrow local min is global min.

Test for convexity.

$$f \text{ is convex} \Leftrightarrow \frac{d^2 f}{d x^2}(x) \geq 0. \quad \forall x.$$

- Sum of convex functions is convex.

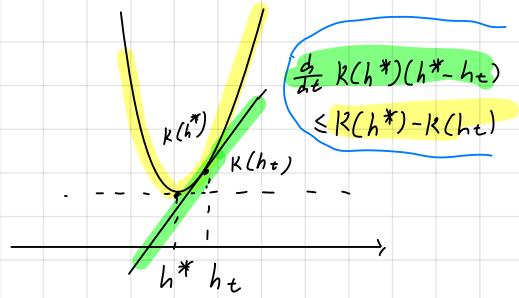
↳ if $L(h, y)$ is convex, then $R(h) = \frac{1}{n} \sum_{i=1}^n L(h, y_i)$ is convex.

- $L_{sq}(h, y) = (y - h)^2$ is convex
- $L_{abs}(h, y) = |y - h|$ is convex

Gradient Descent

$$h_{t+1} = h_t - \alpha \frac{d}{dh} R(h_t).$$

$$0 \leq \frac{d^2}{dh^2} R(h) \leq L \quad \begin{matrix} L=2 & \text{for Huber Risk} \\ \uparrow & \leftarrow \text{concave} \\ \uparrow & \text{smooth} \\ \text{convex} & \end{matrix}$$



$$\begin{aligned} (h_{t+1} - h^*)^2 &= (h_t - \alpha \frac{d}{dh} R(h_t) - h^*)^2 \\ &= (h_t - h^*)^2 - 2\alpha \frac{d}{dh} R(h_t)(h_t - h^*) + \alpha^2 \left(\frac{d}{dh} R(h_t) \right)^2 \\ &\geq R(h^*) - R(h_t) \quad \text{Hv 2 DS} \\ &\leq 2L(R(h_t) - R(h^*)). \end{aligned}$$

$$\leq (h_t - h^*)^2 - 2\alpha(1 - L\alpha)(R(h_t) - R(h^*)).$$

$\alpha = \frac{1}{2L}$ for optimize
 for maximum \Rightarrow less iteration for gradient descent.
 $\frac{1}{2L}$

$$\text{So, let } \alpha = \frac{1}{2L},$$

$$\Rightarrow (h_{t+1} - h^*)^2 \leq (h_t - h^*)^2 - \frac{1}{2L}(R(h_t) - R(h^*)).$$

this gets 0 when $R(h_t) = R(h^*)$, in which case, converges.

$$\Rightarrow \frac{1}{2L}(R(h_t) - R(h^*)) \leq (h_t - h^*)^2 - (h_{t+1} - h^*)^2.$$

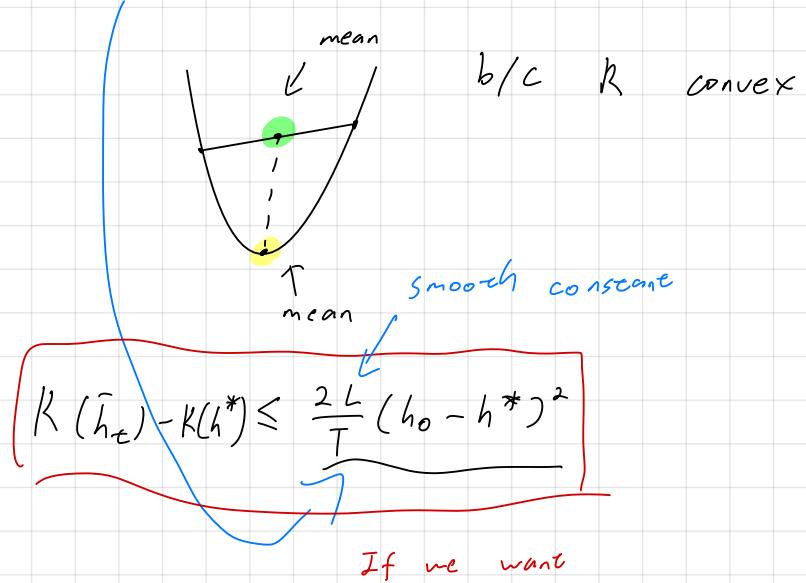
Sum

$$\left(\begin{array}{l} \frac{1}{2L}(R(h_0) - R(h^*)) \leq (h_0 - h^*)^2 - (h_1 - h^*)^2 \\ \frac{1}{2L}(R(h_1) - R(h^*)) \leq (h_1 - h^*)^2 - (h_2 - h^*)^2 \\ \vdots \\ \frac{1}{2L}(R(h_{T-1}) - R(h^*)) \leq (h_{T-1} - h^*)^2 - (h_T - h^*)^2 \end{array} \right) \quad \begin{matrix} \text{telescoping sum} \\ (\text{canceling out terms}). \end{matrix}$$

$$\Rightarrow \frac{1}{2L} \sum_{t=1}^{T-1} (R(h_t) - R(h^*)) \leq (h_0 - h^*)^2 - (h_T - h^*)^2 \leq (h_0 - h^*)^2$$

$$\Rightarrow \frac{1}{2L} \cdot \frac{1}{T} \sum_{t=0}^{T-1} (R(h_t) - R(h^*)) \leq \frac{1}{T} (h_0 - h^*)^2$$

let: $\bar{h}_t = \frac{1}{T} \sum_{t=0}^{T-1} h_t$, then $R(\bar{h}_t) - R(h^*) \leq \frac{1}{T} \sum_{t=0}^{T-1} R(h_t) - R(h^*)$.



Bound $|\bar{x}(D_n^{us}) - \bar{x}(D_\infty^{us})|$

$$R(\bar{h}) - R(h^*) \leq \varepsilon, \quad \text{take } T = \frac{4}{\varepsilon} (h_0 - h^*)^2.$$

$$(\bar{x}(D_n) - \bar{x}(D_\infty))^2 = O(1/n)$$

how large constant

$$(\bar{x}(D_n) - \bar{x}(D_\infty))^2$$

$$\text{Use Var to estimate: } \leq \frac{C}{n} \sum_{i=1}^n (x_i - \bar{x}(D_\infty))^2 \leq \frac{C}{n} \sum_{i=1}^n (x_i - \bar{x}(D_n))^2$$

$$\sum_{i=1}^n (x_i - \bar{x}(D_\infty))^2 \approx \sum_{i=1}^n (x_i - \bar{x}(D_n))^2 = C \cdot \text{Var}(D_n)$$

$$(x_i - \bar{x}(D_\infty) + \bar{x}(D_\infty) - \bar{x}(D_n))^2$$

$$= (x_i - \bar{x}(D_\infty))^2 + (\bar{x}(D_\infty) - \bar{x}(D_n))^2 + 2(\dots)(\dots)$$

$$\frac{1}{n} \cdot (x_i - \bar{x}(D_\infty))^2$$

$$\text{Put together: } (\bar{x}(D_n) - \bar{x}(D_\infty))^2 \leq \frac{C}{n} \cdot \text{Var}(D_n)$$

$$|\bar{x}(D_n) - \bar{x}(D_\infty)| \leq \frac{C}{\sqrt{n}} \cdot \text{STV}(D_n)$$

$$|\bar{x}(D_n) - \bar{x}(D_d)| \leq \frac{\log(1/\delta)}{\sqrt{n}} \cdot \text{STD}(D_n)$$

Sequence Decision Making (Reinforcement Learning)

- Bandits Problem.



Explore then commit

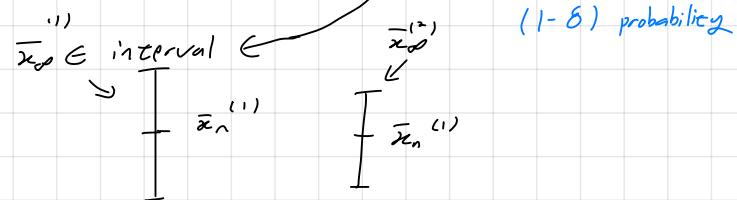
- spend the first N rounds uniformly on the two options.

- Then commit to the one with higher average reward: $\underset{a=1,2}{\operatorname{argmin}} \bar{R}^a$

confidence level
the larger the confidence,
the larger range of error from
true mean

$$W_{\text{ant}}: \Delta \text{reward} = O\left(\frac{\text{STD}(D_N)}{\sqrt{n}}\right) \Leftrightarrow |\bar{x}_n - \bar{x}_{\infty}| \leq \frac{\text{STD}(D_n)}{\sqrt{n}} \cdot \log\left(\frac{1}{\delta}\right)$$

↑
confidence level constant



How to compare \Rightarrow when two interval
do not overlap



one is definitely greater than
the other.

$$\text{by let } \log\left(\frac{1}{\delta}\right) = |\bar{x}_\infty^{(1)} - \bar{x}_\infty^{(2)}| = \Delta \text{reward}$$

↑
 Δ
difference in true mean

How to figure out?

Algorithm:

- keep track N^a

- optimistic estimate: $\hat{R}^a = \bar{R}^a + C \cdot \frac{\text{STD}(D_{N^a})}{N^a}$

- choose action $\underset{a=1,2}{\operatorname{argmin}} \hat{R}^a$

↑
of times
chosen action a

Linear Regression

- Feature - an attribute, information
 - predictor variable - the feature - x
 - response variable - the quantity - y - we are trying to predict.
 - hypothesis function / prediction rule.
- Use loss function to quantify the quality of prediction.
 - Squared loss : $(\text{actual} - \text{predicted})^2$

$$\text{Empirical risk: } R_{\text{sq}}(H) = \frac{1}{n} \sum_{i=1}^n (y_i - H(x_i))^2$$

↓
Smallest mean squared error → good prediction

↑
Not overfitting!

↓
choose H to be certain function!

↓
linear function : $H(x) = w_0 + w_1 x$
 ↑ ↑
 intercept slope

H^* be the linear function that minimizes $R_{\text{sq}}(H) = \frac{1}{n} \sum_{i=1}^n (y_i - H(x_i))^2$ parameter
(least squares regression).
↓ plug in H ↗

Solve for gradient

$$R_{\text{sq}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

$$\frac{\partial R_{\text{sq}}}{\partial w_0} = -\frac{2}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))$$

$$\frac{\partial R_{\text{sq}}}{\partial w_1} = -\frac{2}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i)) x_i$$

Find slope and intercept that minimizes.

Solve for $\frac{\partial R_{sq}}{\partial w_0} = 0$ & $\frac{\partial R_{sq}}{\partial w_1} = 0$

$$w_0 = \frac{1}{n} \sum_{i=1}^n y_i - w_1 \frac{1}{n} \sum_{i=1}^n x_i$$

$$w_0^* = \bar{y} - w_1^* \bar{x}$$

best intercept

$$w_1^* = \frac{\sum_{i=1}^n (y_i - \bar{y}) x_i}{\sum_{i=1}^n (x_i - \bar{x}) x_i}, \text{ where}$$

best slope

$$= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

the process is
called
"fitting to the data".

least square solution
(optimal parameter).

$$H^*(w) = w_0^* + w_1^* x$$

• Correlation coefficient

• a measure of the strength of the linear association of two variables, x & y .

• coefficient, r , the average of the product of x and y , in standard unit

$$r = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

mean
standard deviation.

↓

• r has no unit

• $r = 1 \rightarrow$ perfect positive linear association

• $r = -1 \rightarrow$ perfect negative linear association

• the closer r is to 0 \rightarrow the weaker the linear association

↓ optimal slope for the linear prediction rule w_1^* in r

$$w_1^* = r \frac{s_y}{s_x} \rightarrow \bullet \text{ sign}(r) = \text{sign}(w_1^*)$$

• y value get spread out $\rightarrow s_y \uparrow \rightarrow$ slope \uparrow

• intercept

• x value get spread out $\rightarrow s_x \uparrow \rightarrow$ slope \downarrow

$$H^*(\bar{x}) = \bar{y}$$

Regression

ex. Look at salary (y) as function of years of experience (x).

↑
response variable ↑
predictor variable

$$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

$$y \approx H(x)$$

↑

Hypothesis function / prediction rule.

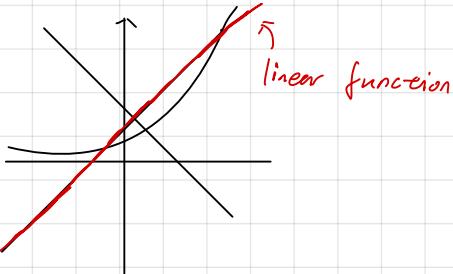
Ex. $y = ax + b$

$$y = e^{ax}$$

$$y = b - x$$

... ...

which $H(x)$ is good?



loss function

$$\text{loss}(H; (x_i, y_i)) = (H(x_i) - y_i)^2 \quad \text{squared loss}$$

$$\begin{aligned} \text{Empirical risk: } R_{\text{sq}}(H, D) &= \frac{1}{n} \sum_{i=1}^n \text{loss}(H; (x_i, y_i)) \\ &= \frac{1}{n} \sum_{i=1}^n (H(x_i) - y_i)^2 \end{aligned}$$

Optimization: choose prediction rule by $\min_{H \in \mathcal{H}} R_{\text{sq}}(H; D)$

One of all loss function: $H(x) = w_1 x + w_0$. Find H^* that minimizes $R_{\text{sq}}(H; D)$.

squared loss ↴



optimize w_1 & w_0 .

$$\min_{w_1, w_0 \in \mathbb{R}} R_{\text{sq}}(H; D)$$

$$= \min \frac{1}{n} \sum_{i=1}^n (H(x_i) - y_i)^2$$

$$= \min \frac{1}{n} \sum_{i=1}^n (w_1 x_i + w_0 - y_i)^2$$

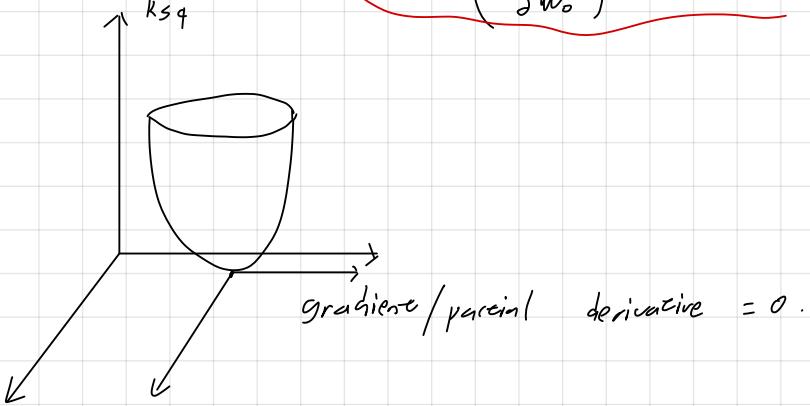
optimal prediction

$$w_1^*, w_0^* = \underset{w_1, w_0}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (w_1 x_i + w_0 - y_i)^2$$

$$\text{let } \nabla_{(w_1, w_0)} \left(\frac{1}{n} \sum_{i=1}^n (w_1 x_i + w_0 - y_i)^2 \right) = 0$$

⇓

$$\nabla_{(w_1, w_0)} \left(f(w_1, w_0) \right) = \begin{pmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_0} \end{pmatrix} = 0 \quad \text{solve for critical point}$$



gradient / parcial derivative = 0.

$$\textcircled{2} \quad \frac{\partial}{\partial w_1} \left(\frac{1}{n} \sum_{i=1}^n (w_1 x_i + w_0 - y_i)^2 \right) = 0.$$

↓

$$\begin{cases} \frac{\partial Ksq((w_1, w_0), v)}{\partial w_0} = 0 & \textcircled{1} \\ \frac{\partial Ksq((w_1, w_0), v)}{\partial w_1} = 0 & \textcircled{2} \end{cases}$$

$$= \frac{1}{n} \sum_{i=1}^n \left(\frac{\partial}{\partial w_1} (w_1 x_i + w_0 - y_i)^2 \right)$$

$$= \frac{1}{n} \sum_{i=1}^n (2 \cdot (w_1 x_i + w_0 - y_i) \cdot x_i)$$

$$= 2 \frac{1}{n} \sum_{i=1}^n (w_1 x_i^2 + w_0 x_i - x_i y_i)$$

$$\textcircled{1} \quad \frac{\partial}{\partial w_0} \left(\frac{1}{n} \sum_{i=1}^n ((w_1 x_i + w_0) - y_i)^2 \right) = 0$$

\textcircled{2} + \textcircled{1} :

$$0 = \frac{1}{n} \sum_{i=1}^n (w_1 x_i^2 + x_i (\bar{y} - w_1 \bar{x}) - x_i y_i)$$

$$0 = \frac{1}{n} \sum_{i=1}^n (w_1 x_i^2 + x_i \bar{y} - x_i w_1 \bar{x} - x_i y_i)$$

$$0 = w_1 \frac{1}{n} \sum_{i=1}^n (x_i^2 - x_i \bar{x}) + \frac{1}{n} \sum_{i=1}^n (x_i \bar{y} - x_i \bar{y})$$

$\bar{x} = x \text{ mean}$

$\bar{y} = y \text{ mean}$

$$w_1 \frac{1}{n} \sum_{i=1}^n (x_i^2 - x_i \bar{x}) = \frac{1}{n} \sum_{i=1}^n (x_i y_i - x_i \bar{y})$$

↓

$$2(w_0 + w_1 \bar{x} - \bar{y}) = 0$$

$w_1 \bar{x} + w_0 = \bar{y}$ optimal solution pass through (\bar{x}, \bar{y}) .

$$w_0 = \bar{y} - w_1 \bar{x}.$$

↓

$$\frac{1}{n} \sum_{i=1}^n x_i (y_i - \bar{y}).$$

↓

$$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Since

$$\frac{1}{n} \sum_{i=1}^n \bar{x} (y_i - \bar{y}) = 0$$

$$\text{Since } \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y}) = 0.$$

↓

Sum of product

↓

Thus, $w_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$

↑
Correlation Variance in x .

Determine the direction of the regression line.
Positive or Negative
since

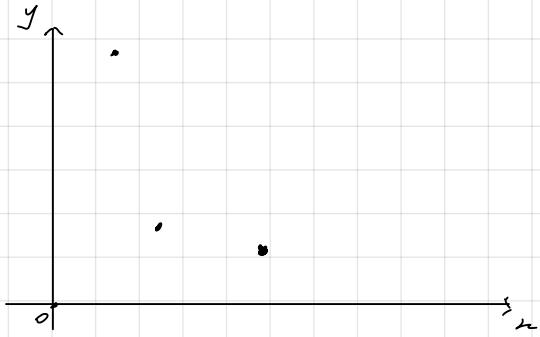
$$\begin{aligned} \sum x_i(y_i - \bar{y}) &= 0 \\ \sum \bar{y}(y_i - \bar{y}) &= 0 \\ &= \sum x_i(y_i - \bar{y}) - \bar{x}y_i - \bar{x}\bar{y} \\ &= \sum y_i(x_i - \bar{x}) = x_i y_i - \bar{x}\bar{y} \\ &= \sum x_i y_i - \bar{x}\bar{y}. \end{aligned}$$

Midterm Review

HW 1

some question in

	x_i	y_i
(x_1, y_1)	3	7
(x_2, y_2)	4	3
(x_3, y_3)	8	2



$$y = H(x) = w_1 x + w_0$$

What is the best prediction rule?

that describes the data set?

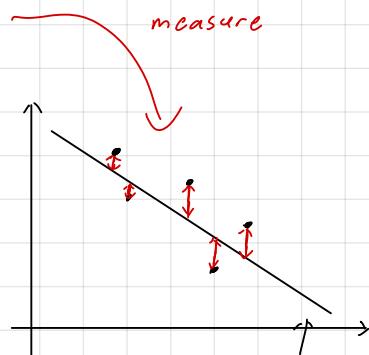
$$\begin{aligned} w_1^* &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} & w_0^* &= \bar{y} - w_1 \cdot \bar{x} \\ &= \frac{(-2)(3) + (-1)(-1) + (3)(-2)}{(-2)^2 + (-1)^2 + (3)^2} & &= 4 - \frac{-11}{14} \cdot 5 \\ &= \frac{-6 + 1 - 6}{4 + 1 + 9} & &= \frac{101}{14} \\ &= \frac{-11}{14} \end{aligned}$$

$$\text{So, } H(x) = -\frac{11}{14}x + \frac{101}{14}.$$

$R_{sq}(H^*; D) = R_{sq}(w_1^*x + w_0^*; D)$. Mean squared error of a data set

$$= \frac{1}{n} \sum_{i=1}^n ((w_1 * x_i + w_0) - y_i)^2$$

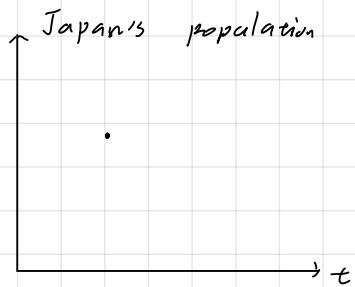
D has two data point



linear prediction rule.

$$R_{sq} = 0$$

ex). Japan's population decrease by half by 2060, will vanish by 2100.



Choices of function close of H.

1. Linear
2. Polynomial
3. Exponential.

next year population

$$X_{t+1} = X_t + d \cdot X_t$$

\uparrow
birth rate

$$X_{t+1} = X_t + d \cdot X_t - \beta X_t$$

$$= (1 + d - \beta) X_t$$

\uparrow death rate
 $= \delta X_t$

$$2. H(x) = w_0 + w_1 x + w_2 x^2 + \dots + w_n x^n$$

$$3. H(x) = c_0 \cdot e^{c_1 x}$$

$$X_t = \delta X_{t-1} = \delta^2 \cdot X_{t-2} = \dots = \delta^t \cdot X_0$$

\uparrow
initial population

$$\underline{X(t) = \delta^t \cdot w}$$

Can we learn (δ, w) from data $\{(t_1, X_1), \dots, (t_n, X_n)\}$?

$$\ln(X(t)) = \ln(\delta^t \cdot w)$$

$$= \ln(\delta^t) + \ln(w)$$

$$= t \cdot \ln(\delta) + \ln(w)$$

$\underbrace{\quad}_{w_1}$ $\underbrace{\quad}_{w_0}$

$$\ln(x(t)) = w_1 t + w_0$$

$s_0,$	t_i	$\ln(x(t))$
,	,	,
,	,	,
,		{}

do least square solution

for w_1 and w_0 .

Conclusion:

$$g(y) = w_0 + w_1 f(x)$$

$$y = f^x \cdot w$$

idea: transform non-linear relationship to linear one. Then do least square solution.

Multi-linear regression

$$y = w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

Want to fit:

$$y = H(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3 + \dots + w_m x^n$$

$$y = H(x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(n)}) \leftarrow \text{multiple feature}$$

Salary $y_0 \in$ location job family, ...

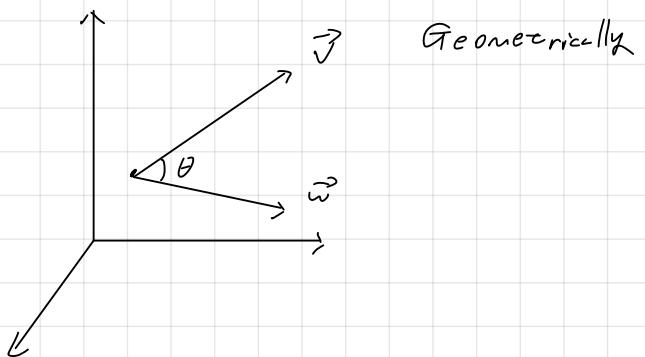
$$= w_0 + w_1 x^{(1)} + \dots + w_n x^{(n)}$$

$$= w_0 + \sum_{i=1}^n w_i \cdot x^{(i)}$$

Linear Algebra (Review):

Vector:

$$\vec{v} = \begin{bmatrix} v_1 \\ \vdots \\ v_d \end{bmatrix} \quad \vec{v} \in \mathbb{R}^d$$



$$\|\vec{v}\|, \frac{\vec{v}^T w}{\|\vec{v}\| \|w\|} = \cos \theta.$$

$$v^T w = \sum_{i=1}^d v_i \cdot w_i = w^T v \leftarrow \text{inner product.}$$

$$[v_1, \dots, v_d] \cdot \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$$

Matrix :

$$A = \begin{pmatrix} \vec{u}_1^T \\ \vdots \\ \vec{u}_n^T \end{pmatrix} = \begin{pmatrix} u_{11} & \dots & u_{1m} \\ \vdots & \ddots & \vdots \\ u_{n1} & \dots & u_{nm} \end{pmatrix} \quad A \in \mathbb{R}^{n \times m}$$

$$= (\vec{v}_1, \dots, \vec{v}_m).$$

Matrix Vector:

$$\vec{y} = A \cdot w = \underbrace{\begin{pmatrix} u_{11} & \dots & u_{1m} \\ \vdots & \ddots & \vdots \\ u_{n1} & \dots & u_{nm} \end{pmatrix}}_{\mathbb{R}^{n \times m}} \underbrace{\begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}}_{\mathbb{R}^m}$$

$$= \begin{bmatrix} u_1^T \cdot w \\ \vdots \\ u_n^T \cdot w \end{bmatrix} \in \mathbb{R}^n$$

$$= \sum_{i=1}^m w_i \cdot \vec{v}_i$$

$$\vec{y} \in \text{span} \{ \vec{v}_1, \dots, \vec{v}_m \}$$

Using L.A. for Multiple linear regression:

$$S_o, \quad w_0 + \sum_{i=1}^d w_i \cdot x^{(i)} \quad \downarrow \quad \vec{w}^T \cdot \vec{x}$$

$$\vec{w} = \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix}, \vec{x} = \begin{bmatrix} x_0 \\ \vdots \\ x_d \end{bmatrix}$$

$$= w_0 x^{(0)} + \sum_{i=1}^d w_i \cdot x^{(i)}, \text{ where } x^{(0)} = 1.$$

$$= \sum_{i=0}^d w_i \cdot x^{(i)}$$

$$= \vec{x}^\top \cdot \vec{w}$$

$$\Leftarrow \text{Thus, } \vec{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}, \vec{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$H(\vec{x}) = \vec{x}^\top \vec{w}$$

$$\vec{x} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix} \in K^{d+1}$$

$$\vec{w} = \begin{pmatrix} w_0 \\ \vdots \\ w_d \end{pmatrix} \in K^{d+1}$$

$$\vec{w}, \vec{x} \in K^{d+1}$$

$$\text{Then, } y = H(\vec{x}) = \vec{x}^\top \cdot \vec{w}$$

A data set will be like

parameter	Salary	YOE	Location	
to learn	y_1	$x_1^{(1)}$	$x_1^{(2)}$	
	y_2	$x_2^{(1)}$	$x_2^{(2)}$	Data point:
	.	.	.	
				(\vec{x}, y) .

$$\text{Loss}(H(\vec{x}, y)) = (y_i - H(\vec{x}_i))^2.$$

$$\text{Loss}(\vec{w}, (\vec{x}_i, \vec{y})) = (y_i - \vec{x}_i^\top \vec{w})^2.$$

$$\text{So, } R(\vec{w}, v) = \frac{1}{n} \sum_{i=1}^n \text{Loss}(\vec{w}, (\vec{x}_i, y_i)).$$

$$= \frac{1}{n} \sum_{i=1}^n (y_i - \vec{x}_i^\top \vec{w})^2$$

$$D = \{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\}$$

$$= \frac{1}{n} \| \vec{v} \|_2^2 \quad \leftarrow \| \vec{v} \|_2^2 = \vec{v}^\top \vec{v} = \sum_{i=1}^n v_i^2$$

let $v_i = y_i - \vec{x}_i^\top \vec{w}$

$$\text{So, } \vec{v} = \begin{bmatrix} y_1 - \vec{x}_1^\top \vec{w} \\ y_2 - \vec{x}_2^\top \vec{w} \\ \vdots \\ y_n - \vec{x}_n^\top \vec{w} \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} \vec{x}_1^\top \vec{w} \\ \vdots \\ \vec{x}_n^\top \vec{w} \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} - \vec{w} \begin{bmatrix} \vec{x}_1 \\ \vdots \\ \vec{x}_n \end{bmatrix}$$

$$\vec{X} \in K^{n \times (d+1)}$$

$$= \vec{y} - \vec{X} \cdot \vec{w}$$

$\vec{y} \in K^n$

$\vec{X} \in K^{n \times (d+1)}$

Design Matrix

Thus,

$$y = H(\vec{w}) = \vec{x}^T \cdot \vec{w}$$

$$R(\vec{w}; D) = \frac{1}{n} \| \vec{y} - \vec{x} \cdot \vec{w} \|^2$$

$$= \frac{1}{n} (\vec{y} - \vec{x} \cdot \vec{w})^T (\vec{y} - \vec{x} \cdot \vec{w}).$$

$$= \frac{1}{n} (\vec{y}^T - \vec{w}^T \vec{x}^T) (\vec{y} - \vec{x} \cdot \vec{w}).$$

$$= \frac{1}{n} (\vec{y}^T \vec{y} - \underbrace{\vec{y}^T \vec{x} \vec{w}}_{\text{scalar}} - \underbrace{\vec{w}^T \vec{x}^T \vec{y}}_{\text{scalar}} + \vec{w}^T \vec{x}^T \vec{x} \vec{w})$$

$$\vec{w} = \arg \min_{\vec{w} \in R^{d+1}} R(\vec{w}, D)$$

$$\vec{y}^T \vec{x} \vec{w} = (\vec{y}^T \vec{x} \vec{w})^T = \vec{w}^T \vec{x}^T \vec{y}$$

$$= \frac{1}{n} (\vec{y}^T \vec{y} - 2 \vec{w}^T \vec{x}^T \vec{y} + \vec{w}^T \vec{x}^T \vec{x} \vec{w}).$$

$$S_o, \nabla_{\vec{w}} R_{\text{sq}}(\vec{w}, D) \in R^d$$

$$= \frac{1}{n} (\nabla_{\vec{w}} \vec{y}^T \vec{y} - 2 \nabla_{\vec{w}} \vec{w}^T \vec{x}^T \vec{y} + \nabla_{\vec{w}} \vec{w}^T \vec{x}^T \vec{x} \vec{w}).$$

$$= X^T \vec{y} \in R^{d+1}$$

$$= 2 X^T \vec{x} \vec{w}$$

$$\text{Since } \nabla_{\vec{w}} \vec{w}^T \vec{b}$$

$$= \begin{pmatrix} \frac{\partial}{\partial w_0} \\ \vdots \\ \frac{\partial}{\partial w_d} \end{pmatrix} (w_0 b_0 + \dots + w_d b_d)$$

$$= \begin{pmatrix} b_0 \\ \vdots \\ b_d \end{pmatrix} = \vec{b}$$

$$\text{Thus, } 0 = \frac{1}{n} (-2 X^T \vec{y} + 2 X^T \vec{x} \vec{w})$$

$$0 = -X^T \vec{y} + X^T \vec{x} \vec{w}$$

$$\underline{X^T \vec{y} = X^T \vec{x} \vec{w}}$$

Normal Equation

$$\begin{matrix} X^T X & \xrightarrow{\text{invertible}} & \\ \begin{matrix} \uparrow \\ R^{(d+1) \times n} \end{matrix} & & \begin{matrix} \uparrow \\ R^{n \times (d+1)} \end{matrix} \\ \downarrow & & \end{matrix}$$

$$R^{(d+1) \times (d+1)}$$

Square matrix R full rank

\uparrow

(columns of $X^T X$ linearly independent)

Empirical Risk

$$\vec{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix}$$

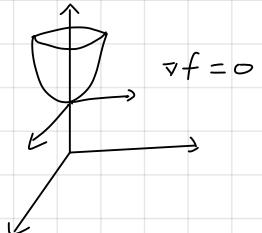
$$= \begin{pmatrix} x_{1(1)} & \cdots & x_{1(d)} \\ \vdots & \ddots & \vdots \\ \vdots & \ddots & \vdots \\ x_{n(1)} & \cdots & x_{n(d)} \end{pmatrix}$$

Design
Matrix

From one-dimensional LS

$$\begin{cases} \frac{\partial}{\partial w_0} R(w_0, w_1, D) = 0 \\ \frac{\partial}{\partial w_1} R(w_0, w_1, D) = 0 \end{cases}$$

multi-dimensional
w.r.t vector



$$\begin{pmatrix} \frac{\partial}{\partial w_0} R(\vec{w}, D) \\ \vdots \\ \frac{\partial}{\partial w_d} R(\vec{w}, D) \end{pmatrix} = 0$$

$$\begin{pmatrix} \frac{\partial}{\partial w_0} \\ \vdots \\ \frac{\partial}{\partial w_d} \end{pmatrix} R(\vec{w}, D) = 0$$

$$\nabla_{\vec{w}} R(\vec{w}, D) = 0$$

$$\begin{matrix} \uparrow \\ P \\ \in R^{d+1} \end{matrix} \quad \begin{matrix} \uparrow \\ R \end{matrix}$$

thus, $\exists (X^T X)^{-1}$,

$$(X^T X)^{-1} X^T X = I$$

Our normal equation becomes

$$\vec{w} = (X^T X)^{-1} X^T \vec{y}$$

why matter?

$$(X^T X) \vec{w} = \sum_{i=0}^d w_i \cdot b_i$$

if dependent, $\exists \vec{w}' \neq 0$, s.t. $(X^T X) \vec{w}' = 0$.

linear combination of $X^T X$
s.t. $= 0$.

Q: When is $X^T X$ full rank?

① X is square matrix, $n = d+1$

$$X^T X = \begin{pmatrix} 1 & \dots & 1 \\ x_{1(1)} & \dots & x_{n(1)} \\ \vdots & \vdots & \vdots \\ x_{1(d)} & \dots & x_{n(d)} \end{pmatrix} \begin{pmatrix} 1 & x_{1(1)} & \dots & x_{n(1)} \\ \vdots & \vdots & & \vdots \\ 1 & x_{1(d)} & \dots & x_{n(d)} \end{pmatrix}$$

$$X^T(X\vec{w}) \rightarrow X^T \text{ full rank} \rightarrow X \text{ full rank}$$

(data point not repeating
or not linearly combination of other data points)

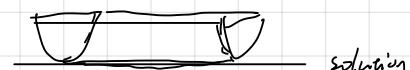
② data point linearly dependent

⇓

$$n < d+1$$

over-parametrized

convex



⇓

infinity many solutions.

always solution to
normal equation

Meaning of optimal solution

$$\vec{y} \approx H(\vec{x}) = \vec{X}^T \vec{w}^*$$

↳ average salary of industry

$$\vec{y} = w_0^* + w_1^* x_{1(1)} + w_2^* x_{1(2)} + w_3^* x_{1(3)} + \dots + w_d^* x_{1(n)}$$

Salary

↑
Y.O.E.

↑
years of between job

Standardize (z-score)

$$w_3^* : 100,000 \rightarrow 100,001$$

$$w_1^* : 1 \text{ year} \rightarrow 2 \text{ years}$$

↓

$$w_1^* \gg w_3^*$$

$$w_3^* : K\$100 \rightarrow K\$101$$

scalar

If we standardize)

w^* reflect correlation.

ex) $H(x) = w_0 + w_1 x + w_2 x^2 + w_3 x^3 + \dots + w_d x^d$ ← linear regression okay

limitation of linear regression?

No, b/c Taylor series idea → polynomial approximates any function.

$$H(n) = w_0 + w_1 \sin(n\pi) + d_1 \cos(n\pi) + w_2 \sin(2n\pi) + d_2 \cos(2n\pi) + \dots$$

↑
can be used to model periodical function.

Problem:

- Overfitting
- efficiency

ex). Free falling obj.

$$H(n) = w_0 + w_1 n + w_2 n^2 + \dots -$$

not necessary

$$H(\vec{x}) = w_0 + w_1^T \vec{x} + \vec{x}^T W_2 \vec{x}$$

hard to solve (high dimensional)

Hierarchical model (deep learning model) classification problem

$h^t(\vec{x}) = \sigma(W^t \vec{x})$

$W^t \in \mathbb{R}^{d \times d}$

depends on how large the next layer you want

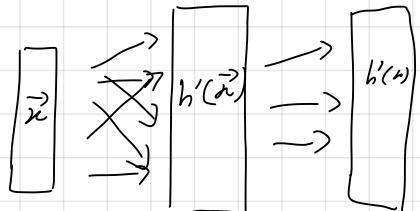
$$h^2(\vec{x}) = \sigma(w^2 h^1(\vec{x}))$$

:

$$h^t(\vec{x}) = \sigma(w^t h^{t-1}(\vec{x}))$$

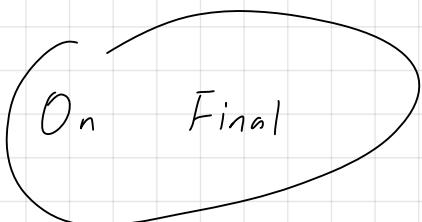
,

$$h^T(\vec{x}) = \underbrace{\vec{w}^T}_{\equiv} \underbrace{h^{T-1}(\vec{x})}_{\equiv}$$

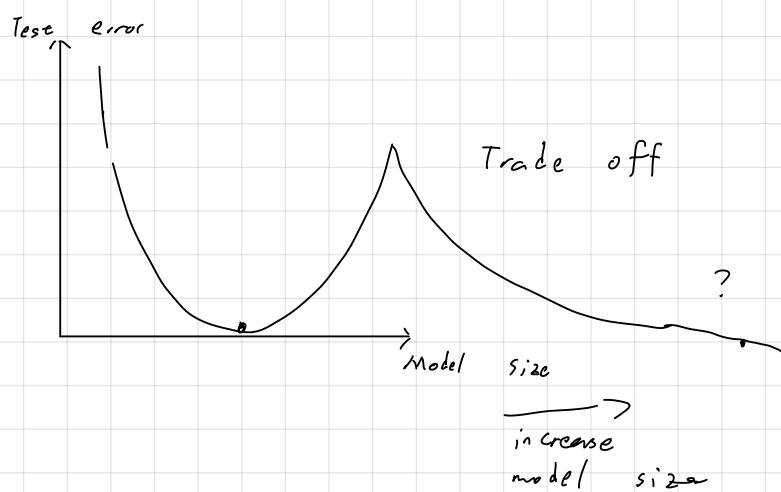


$$\delta_i(\vec{z}) = \begin{cases} 0, & z_i < 0 \\ z_i, & z_i \geq 0 \end{cases}$$

Rectified Linear Unit

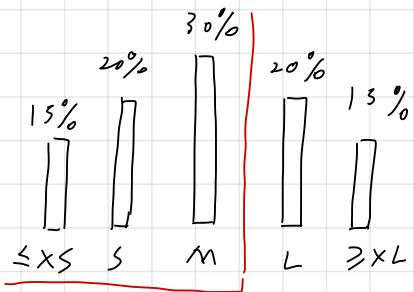


Vector Calculus



Probability theory

distribution are similar



Probability Statement:

Outcome (" $\leq xS$ ", "S", "M", "L", " $\geq xL$ ")

Sample space $S = \{ \quad \}$

Event E is a subset of sample space

Probability P : a function that maps event to a real number.

$$\textcircled{1} \quad 0 \leq P(E) \leq 1, \forall E \subseteq S$$

$$\textcircled{2} \quad P(E_1 \cup E_2) = P(E_1) + P(E_2), \forall E_1 \text{ and } E_2 \text{ not intersecting each other.}$$

$$\forall E_1 \cap E_2 = \emptyset$$

$$E_1 = \{ " \leq xS ", "S", "M" \}$$

$$E_2 = \{ "M", "L" \}$$

$$E_1 \cup E_2 = \{ " \leq xS ", "S", "M", "L" \}$$

$$E_1 \cap E_2 = \{ "M" \}$$

point :



Discrete Sample Space:

$$P(\bigcup_{i=1}^{n-1} \{S_i\}) = P(\bigcup_{i=1}^{n-1} \{S_i\} \cup S_n)$$

$$= P(\bigcup_{i=1}^{n-1} \{S_i\}) + P(S_n)$$

$$= \sum_{i=1}^n P(\{S_i\})$$

$$= 1$$

interval :

\Rightarrow Probability P : a function maps outcomes to real number

$$\textcircled{1} \quad 0 \leq P(\{S_i\}) \leq 1, \forall S_i \in S.$$

$$\textcircled{2} \quad \sum_{i=1}^n P(\{S_i\}) = 1$$

$$\text{ex). } E_1 = \{1, 2, 3\}, E_2 = \{1, 5, 6\}$$

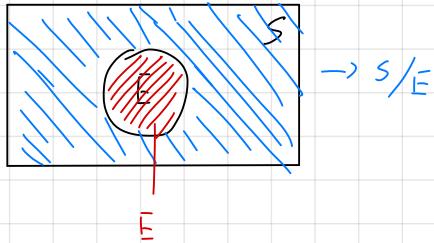
$$\underline{P}(E_1) = \frac{1}{2} = \underline{P}(E_2) \rightarrow \text{Uniform Probability}$$

Venn Diagram

Event E . $\tilde{E} = S|_E$, $\underline{P}(\tilde{E}) = 1 - \underline{P}(E)$.

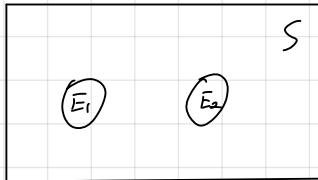
complement of E

$$\underline{P}(E) + \underline{P}(\tilde{E}) = \underline{P}(E \cup \tilde{E}) = \underline{P}(S) = 1.$$



① $\underline{P}(E_1 \cup E_2) = \underline{P}(E_1) + \underline{P}(E_2)$, if $E_1 \cap E_2 = \emptyset$

↑
mutually exclusive.



② $\underline{P}(\tilde{E}_1 \cup \tilde{E}_2) = \underline{P}(E_1) + \underline{P}(E_2) - \underline{P}(E_1 \cap E_2)$.

$$\underline{P}(E_1 \cup E_2) = \underline{P}(E_1 \setminus E_1 \cap E_2) + \underline{P}(E_2 \setminus E_1 \cap E_2) + \underline{P}(E_1 \cap E_2).$$

$$= \underline{P}(E_1) - \underline{P}(E_1 \cap E_2) + \underline{P}(E_2) - \underline{P}(E_1 \cap E_2) + \underline{P}(E_1 \cap E_2)$$

$$= \underline{P}(E_1) + \underline{P}(E_2) - \underline{P}(E_1 \cap E_2).$$

ex). 0.3 prob. roommates home $\rightarrow A$

0.2 prob. roommates partner home $\rightarrow B$

0.1 prob. both of them home $\rightarrow A \cap B$

Prob no one is at home \leftarrow let C : no one is at home

$$\underline{P}(C) = 1 - \underline{P}(\tilde{C}) = 1 - \underline{P}(A \cup B)$$

$$= 1 - (\underline{P}(A) + \underline{P}(B) - \underline{P}(A \cap B))$$

$$= 1 - 0.3 - 0.2 + 0.1 \\ = 0.6.$$

Sample space = $\{(R \text{ home}, P \text{ home}), (R \text{ not h}, P h), (R h, P \text{ not h}), (R \text{ not h}, P \text{ not h})\}$

③ $\underline{\underline{P}}(E_1 \cap E_2) = \underline{\underline{P}}(E_1) \underline{\underline{P}}(E_2 | E_1) = \underline{\underline{P}}(E_2) \cdot \underline{\underline{P}}(E_2 | E_1)$

↑
prob of E_2 , given condition on E_1 .

$$E_1, E_2 \rightarrow \underline{\underline{P}}(E_1 \cap E_2) = \underline{\underline{P}}(E_1) \underline{\underline{P}}(E_2)$$

$\stackrel{\text{independent}}{\Downarrow}$
 $\underline{\underline{P}}(E_2) = \underline{\underline{P}}(E_2 | E_1).$

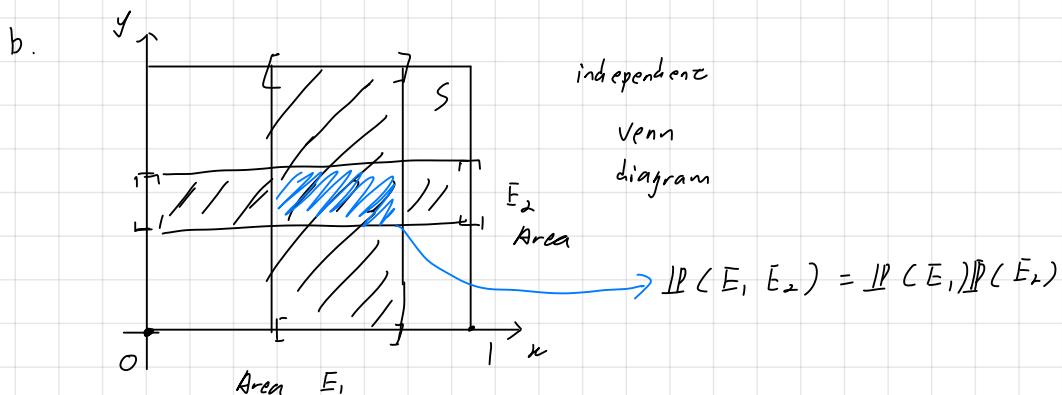
④ Are mutually exclusive event independent?

$$\underline{\underline{P}}(E_1 \cap E_2) = \underline{\underline{P}}(\emptyset) = 0$$

$$\underline{\underline{P}}(E_1 | E_2) = 0 \quad \text{if} \quad \underline{\underline{P}}(E_2) \neq 0.$$

$$= \frac{\underline{\underline{P}}(E_1 \cap E_2)}{\underline{\underline{P}}(E_2)} \leftarrow = 0 \quad \leftarrow \neq 0$$

a. if $\underline{\underline{P}}(E_1) \neq 0$ & $\underline{\underline{P}}(E_2) \neq 0$, $E_1 \cap E_2 = \emptyset \Rightarrow E_1, E_2$ not independent



ex) Roll a die 3 times

prob $\{$ face 1 never appears $\}$

prob $\{$ face 1 appear at least one $\} \leftarrow B_1$

sample space for each roll : $\{1, \dots, 6\}$

$$A_1 = \{ \text{roll 1 = face 1} \} . \quad \underline{\mathbb{P}}(A_1) = \frac{1}{6}$$

$$\bar{A}_1 = \{ \text{roll 1} \neq \text{face 1} \} . \quad \underline{\mathbb{P}}(\bar{A}_1) = \frac{5}{6}$$

$$\underline{\mathbb{P}}(\bar{A}_1 \cap \bar{A}_2 \cap \bar{A}_3) = \underline{\mathbb{P}}(\bar{A}_1) \cdot \underline{\mathbb{P}}(\bar{A}_2) \cdot \underline{\mathbb{P}}(\bar{A}_3)$$

$$= \left(\frac{5}{6}\right)^3$$

$$\underline{\mathbb{P}}(\bar{B}_1) = 1 - \underline{\mathbb{P}}(\bar{A}_1 \cap \bar{A}_2 \cap \bar{A}_3)$$

$$= 1 - \left(\frac{5}{6}\right)^3$$

ex) Two pets Equal prob. $\{$ cat, dog $\}$.

Q1 Prob both are dogs, given oldest is dog? (D, D)

Q2 ----, at least one of them is dog? (D, C), (C, D)
weaker condition (D, D).

Q1 Event $B = \{ \text{both are dogs} \} \quad \underline{\mathbb{P}}(B \cap A_1) = \underline{\mathbb{P}}(B) = \left(\frac{1}{2}\right)^2 = \frac{1}{4}$

Event $A_1 = \{ \text{oldest is dog} \} \quad \underline{\mathbb{P}}(A_1) = \frac{1}{2}$

$$\underline{\mathbb{P}}(B|A_1) = \frac{\underline{\mathbb{P}}(B \cap A_1)}{\underline{\mathbb{P}}(A_1)} = \frac{\frac{1}{4}}{\frac{1}{2}} = \frac{1}{2}$$

Q2 Event $A_2 = \{ \text{at least one is dog} \} \quad \underline{\mathbb{P}}(A) = 1 - \underline{\mathbb{P}}(\bar{A}_2) = \frac{3}{4}$

$\bar{A}_2 = \{ \text{none of them are dogs} \} \quad \underline{\mathbb{P}}(\bar{A}_2) = \frac{1}{4}$

$$\underline{\mathbb{P}}(B|A_2) = \frac{\underline{\mathbb{P}}(B \cap A_2)}{\underline{\mathbb{P}}(A_2)} = \frac{\frac{1}{4}}{\frac{3}{4}} = \frac{1}{3}$$

$$\underline{\mathbb{P}}(B \cap A_2) = \underline{\mathbb{P}}(B) = \frac{1}{4}$$

Condition Probability

$$\underline{P(B|A)}$$

↑
condition

Law of Total Probability

ex)

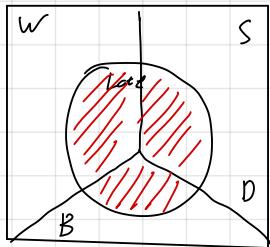
Getting to campus

	Late	on time
walk :	0.06	0.24
bike :	0.03	0.07
drive :	0.36	0.24

Q1 Prob a random is late?

$$P(\text{Late}) = P(\text{Late} \cap \text{walk}) + P(\text{Late} \cap \text{bike}) + P(\text{Late} \cap \text{drive}).$$

$$= P(\text{Late} | \text{walk}) P(\text{walk}) + P(\text{Late} | \text{bike}) P(\text{bike}) + P(\text{Late} | \text{drive}) P(\text{drive}).$$



$$(L \cap D) \cup (L \cap B) \cup (L \cap W) = L \quad \leftarrow \text{disjoint}$$

Property $(A \cap B) \cup (A \cap C) = A \cap (B \cup C)$

$$P(L \cap D) + P(L \cap B) + P(L \cap W) = P(L)$$

Partition $\{E_1, \dots, E_k\}$ is partition of S if

$$\textcircled{1} \quad P(E_i \cap E_j) = 0, \quad \forall i \neq j. \quad (E_i \cap E_j = \emptyset, \quad \forall i \neq j).$$

$$\textcircled{2} \quad \sum_{i=1}^k P(E_i) = 1 \quad \left(\bigcup_{i=1}^k E_i = S \right).$$

$$\textcircled{3} \quad P(A) = \sum_{i=1}^k P(A \cap E_i).$$

$$= \sum_{i=1}^k P(A | E_i) P(E_i)$$

$$Q2 \quad \underline{P(Walk | Late)} = \frac{\underline{P(Late \cap Walk)}}{\underline{P(Late)}} = \frac{0.06}{0.06 + 0.03 + 0.36}$$

$$Q3 \quad \underline{P(Late)} = 0.45$$

$$\underline{P(Walk)} = 0.3$$

$$\underline{P(Late | Walk)} = 0.2$$

$$\underline{P(Walk | Late)} = \frac{\underline{P(Walk \cap Late)}}{\underline{P(Late)}} = \frac{\underline{P(Late | Walk)} \underline{P(Walk)}}{\underline{P(Late)}}$$



Bayes Theorem

Statistical Model

$$\underline{P(A | B)} = \frac{\underline{P(B | A)} \underline{P(A)}}{\underline{P(B)}}$$

$$- \frac{\underline{P(Data | Parameter)} \underline{P(Parameter)}}{\underline{P(Data)}} = \underline{P(Parameter | Data)}$$

↑
prior belief $\underline{P(Parameter | Data)}$

ex) Product detects Covid-19 w/ 95% prob.

15% prob. turning positive for people w/o covid-19.

10% have covid-19

Q: Prob. people get covid-19, if they test positive?

T: {test positive}

H: {have covid-19}

$$\underline{P(T)} = \underline{P(T | H)} \underline{P(H)} + \underline{P(T | H^c)} \underline{P(H^c)}$$

$$\underline{P(T | H)} = 0.95$$

$$\underline{P(T | H^c)} = 0.15$$

$$\underline{P(H)} = 0.1 \quad 0.95 \cdot 0.1$$

$$\underline{P(H | T)} = \frac{\underline{P(T | H)} \underline{P(H)}}{\underline{P(T)}}$$

$$\left. \right) = 0.95 \cdot 0.1 + 0.15 \cdot 0.9$$

$$\left. \right) \approx 0.4$$

Differential Privacy

338 website takes a poll

Republ. Dem.

49% 51%

- Differential attack.

- Privacy - accuracy trade off

- Guard against privacy leakage

Flip a coin — "head" → answer faithfully

— "tail" → answer 50% - 50%

Republ. Dem.

Q: How safe is this?

Denote "Rep" - 0

"Dem" - 1

Private preference - $X_n \in \{0, 1\}$ private

Observation - $Y_n \in \{0, 1\}$ ✓

Question: $\Pr(X_n = 1 | Y_n = 0) = ? \longrightarrow \Pr(\text{prefer Dem.} | \text{Answered "Rep"})$.

$$= \frac{\Pr(Y_n = 0 | X_n = 1) \Pr(X_n = 1)}{\Pr(Y_n = 0)} \quad \begin{matrix} \uparrow \\ \text{observed} \end{matrix}$$

$$\Pr(Y_n = 0 | X_n = 1) = \Pr(\text{faithful}) \Pr(Y_n = 0 | X_n = 1, \text{faithful}) + \Pr(\text{not faithful}) \Pr(Y_n = 0 | X_n = 1, \text{not faithful})$$

$$= \frac{1}{2} \cdot 0 + \frac{1}{2} \cdot \frac{1}{2} \quad \Pr(\text{Not faithful}) \Pr(Y_n = 0 | X_n = 1, \text{not faithful}) = \frac{3}{4}$$

$$= \frac{1}{4}$$

Law of Total Probability

$$\Pr(X_n = 1) = 0.51 \quad \left. \begin{array}{l} \text{Assume} \\ \text{Distribution} \\ \text{of vote} \end{array} \right| \quad \Pr(Y_n = 0) = \Pr(Y_n = 0 | X_n = 0) \Pr(X_n = 0) + \Pr(Y_n = 0 | X_n = 1) \Pr(X_n = 1)$$

$$= \frac{3}{4} (0.49) + \frac{1}{4} (0.51)$$

$$= 0.495$$

$$\frac{\underline{P}(Y_1 = 0 \mid X_1 = 1) \underline{P}(X_1 = 1)}{\underline{P}(Y_1 = 0)} = \frac{\frac{1}{4} - 0.51}{0.495} \approx \frac{1}{4}.$$

$$\underline{P}(X_1 = 0 \mid Y_1 = 0) = 1 - \underline{P}(X_1 = 1 \mid Y_1 = 0) \approx \frac{3}{4}. \uparrow$$

prob that attacker success.

Q: Can we recover the poll result?

Don't know $\underline{P}(X=1)$ but know $\underline{P}(Y=1)$.

$$\underline{P}(A \mid B) \stackrel{?}{=} \underline{P}(\bar{A} \mid B)$$

Expectation:

$$\frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} \left(\sum_{x_i \in D_0} x_i + \sum_{x_j \in D_1} x_j \right)$$

$$n = n_0 + n_1 = \frac{1}{n} (n_0 \cdot 0 + n_1 \cdot 1)$$

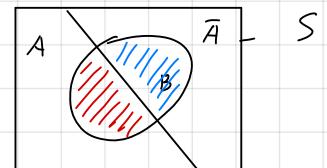
$$\begin{matrix} \uparrow & \uparrow \\ \text{Rep. } D_0 & D_1 \end{matrix} = \frac{n_0}{n} \cdot 0 + \frac{n_1}{n} \cdot 1$$

$$\text{As } n \rightarrow \infty = \underline{P}(X = 0) \cdot 0 + \underline{P}(X = 1) \cdot 1$$

$$= \underline{P}(X = 1)$$

$$\text{So, } \frac{1}{n} \sum_{i=1}^n Z_i \xrightarrow{n \rightarrow \infty} \sum_{k=1}^K \underline{P}(Z = a_k) a_k$$

\uparrow pmf
Expectation. \uparrow value
pmf take



$$(A \cap B) \cup (\bar{A} \cap B) = B$$

$$\begin{aligned} & \uparrow \\ & \underline{P}(A \cap B) + \underline{P}(\bar{A} \cap B) \\ & = \underline{P}(B) \end{aligned}$$

$$\text{Similarly, } \frac{1}{n} \sum_{i=1}^n Y_i = \underline{P}(Y = 0) \cdot 0 + \underline{P}(Y = 1) \cdot 1$$

$$\uparrow = \underline{P}(Y = 1) \quad \leftarrow \text{law of total probability}$$

Expectation of Y

$$= \underline{P}(Y = 1 \mid X = 0) \underline{P}(X = 0) + \underline{P}(Y = 1 \mid X = 1) \underline{P}(X = 1)$$

$$= \frac{1}{4} \underline{P}(X = 0) + \frac{3}{4} \underline{P}(X = 1).$$

$$= \frac{1}{4} (1 - \underline{P}(X = 1)) + \frac{3}{4} \underline{P}(X = 1)$$

$$= \frac{1}{4} - \frac{1}{4} \underline{P}(X = 1) + \frac{3}{4} \underline{P}(X = 1)$$

$$= \frac{1}{4} + \frac{1}{2} \mathbb{P}(X=1).$$

\uparrow
Expectation of X

$$\mathbb{P}(Y=1) = \frac{1}{4} + \frac{1}{2} \mathbb{P}(X=1)$$

This is how to recover the result!

$$\mathbb{P}(X=1) = 2(\mathbb{P}(Y=1) - \frac{1}{4}).$$

Probability on continuous space

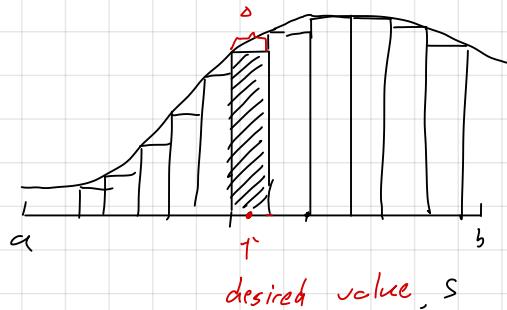
Discrete: $s \in S$, ① $0 \leq \mathbb{P}(s) \leq 1$.

$$\textcircled{2} \quad \sum_{s \in S} \mathbb{P}(s) = 1.$$

$$\lim_{k \rightarrow \infty} \sum_{\Delta \rightarrow 0}^k \mathbb{P}\left(X \in [s_k - \frac{\Delta}{2}, s_k + \frac{\Delta}{2}]\right) = 1.$$

$$= \lim_{\Delta \rightarrow 0} \sum_{k=1}^K p(s_k) \cdot \Delta$$

$$= \int_a^b p(s) ds.$$



$$\lim_{\Delta \rightarrow 0} \mathbb{P}\left(X \in [s - \frac{\Delta}{2}, s + \frac{\Delta}{2}]\right) = p(s)$$



Probability Density Function

Continuous:

$$\textcircled{1} \quad p(s) \geq 0$$

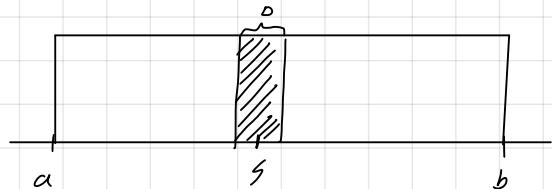
$$\textcircled{2} \quad \int_{S \in S} p(s) ds = 1.$$

Uniform probability case

* Uniform Distribution.



$$\text{p.d.f. } p(s) = \frac{1}{b-a}$$



$$\mathbb{P}\left(X \in [s - \frac{\Delta}{2}, s + \frac{\Delta}{2}]\right) = \frac{\Delta}{b-a}$$

• Normal Distribution.

p.d.f. defined on \mathbb{R} .

$$\text{p.d.f. } p(s) = \frac{1}{\sqrt{2\pi} \cdot \sigma} \exp\left(-\frac{(s-\mu)^2}{2\sigma^2}\right)$$

$$\int p(s) ds = 1$$

$$\begin{aligned} &= \int_{\mathbb{R}} \exp\left(-\frac{(s-\mu)^2}{2\sigma^2}\right) ds \\ &= \sqrt{2\pi} \cdot \sigma \end{aligned}$$

μ : mean / Expectation of distribution



$$\textcircled{1} \text{ Discrete Expectation : } \frac{1}{n} \sum_{i=1}^n x_i \xrightarrow{n \rightarrow +\infty} E[X]$$

$$E[X] = \sum_{k=1}^K P(X = s_k) \cdot s_k$$

$$\begin{aligned} \textcircled{2} \text{ Continuous Expectation: } E[X] &= \lim_{\Delta \rightarrow 0} \sum_{k=1}^K P(X \in [s_k - \frac{\Delta}{2}, s_k + \frac{\Delta}{2}]) \cdot s_k \\ &= \lim_{\Delta \rightarrow 0} \sum_{k=1}^K p(s_k) \Delta \cdot s_k \\ &= \int_{-\infty}^{\infty} p(s) \cdot s \, ds. \end{aligned}$$

$$\text{For normal distribution: } E[X] = \int_{-\infty}^{\infty} p_N(s) \cdot s \, ds = \mu$$

σ : STD, σ^2 : Variance

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\begin{aligned} E[f(x)] \\ = \sum_{k=1}^K P(X = s_k) \cdot f(s_k). \end{aligned}$$

$$\begin{aligned} \int_{-\infty}^{\infty} p(s) \cdot (s - \mu)^2 \, ds &= E[(X - \mu)^2]. \\ \downarrow n \rightarrow \infty \\ \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 &\xrightarrow{n \rightarrow \infty} \sum_{k=1}^K P(X = s_k) \cdot (s_k - \mu)^2. \end{aligned}$$

$$\frac{1}{\sqrt{2\pi} \cdot 6} \int_{-\infty}^{\infty} \exp(-\frac{(s-\mu)^2}{2\sigma^2}) (s-\mu)^2 \, ds.$$

$$\text{Let } y = s - \mu$$

$$\frac{1}{\sqrt{2\pi} \cdot 6} \int_{-\infty}^{\infty} \exp\left(-\frac{y^2}{2\sigma^2}\right) y^2 \, dy.$$

do integral by parts twice

$$\exp(-y^2/2\sigma^2) = \exp(-y^2/2\sigma^2) \cdot \left(\frac{-2y}{2\sigma^2}\right) \, dy$$

$$\begin{aligned}
 & \frac{1}{\sqrt{2\pi} \cdot 6} \int_{-\infty}^{\sigma} \exp\left(\frac{-y^2}{2 \cdot 6^2}\right) (y^2) dy \\
 &= \frac{1}{\sqrt{2\pi} \cdot 6} \int_{-\infty}^{\sigma} \left(-\frac{6^2}{y}\right) \cdot d\exp\left(\frac{-y^2}{2 \cdot 6^2}\right) \\
 &= -\frac{1}{\sqrt{2\pi} \cdot 6} \int_{-\infty}^{\sigma} \exp\left(\frac{-y^2}{6^2}\right) \cdot (-6^2) dy \\
 &= 6^2 \cdot \frac{1}{\sqrt{2\pi} \cdot 6} \int_{-\infty}^{\sigma} \exp\left(-\frac{y^2}{6^2}\right) dy \\
 &= 6^2 \cdot \underbrace{\frac{1}{\sqrt{2\pi} \cdot 6}}
 \end{aligned}$$

$$x \sim N(\mu, \sigma^2) \quad p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu)^2}{2\sigma^2})$$

$$E[X] = \int_{-\infty}^{\infty} p(x) \cdot x \, dx = \mu$$

$$\text{Var}(x) = E[(x - \mu)^2] = \int_{-\infty}^{\infty} p(x) (x - \mu)^2 dx = \sigma^2.$$

$$E[f(x)] = \int_{-\infty}^{\infty} p(x) f(x) dx$$

$$\int_{-\infty}^{\infty} p(x, y) \, dy = p(x)$$

$$E[X+Y] = \iint_K p_{(X,Y)}(x,y) \cdot (x+y) dx dy$$

$$P_{x,y}(x,y) = \lim_{\Delta \rightarrow 0} \frac{\mathbb{P}\left(X \in [x - \frac{\Delta}{2}, x + \frac{\Delta}{2}] \wedge Y \in [y - \frac{\Delta}{2}, y + \frac{\Delta}{2}]\right)}{\Delta^2}$$

$$= \iint_{\mathbb{R}^2} p_{(X,Y)}(x,y) \cdot x \, dx dy + 1$$

$\left| \int_{-\infty}^{\infty} p_{X,Y}(x,y) dy \right| = \sum_{n=1}^{\infty} \underbrace{\Pr(X \in [x - \frac{a}{2}, x + \frac{a}{2}] \cap Y \in [y - \frac{a}{2}, y + \frac{a}{2}])}_{\geq 2}$

$$\iint_{\mathbb{R}^2} p_{(X,Y)}(x,y) \, dy dx$$

$$= \int_K x \int_K p_{X,Y}(x,y) dy dn +$$

$$= \frac{\pi (n \cdot l \cdot 2, \dots, 2, 1, \dots)}{0}$$

$$= \frac{\pi (n \cdot L^2 - 2, \dots, 2, 1, \dots)}{D}$$

$$= \frac{\mathbb{I}(X \in [x - \frac{\delta}{2}, x + \frac{\delta}{2}])}{\delta} = p_X(x)$$

$$= \int_K x p_x(x) dx + \int_K y p_y(y) dy$$

$$= \frac{\mathbb{I}(X \in [x - \frac{\delta}{2}, x + \frac{\delta}{2}])}{\delta} = p_X(x)$$

Independence & Expectation

$$\mathbb{P}(A \cap B) = \mathbb{P}(A) \mathbb{P}(B)$$

$$P_{(X,Y)}(x,y) = P_X(x) P_Y(y)$$

b/c

$$P_{X,Y}(x,y) = \lim_{\Delta \rightarrow 0} \frac{\mathbb{P}(X \in [x - \frac{\Delta}{2}, x + \frac{\Delta}{2}], Y \in [y - \frac{\Delta}{2}, y + \frac{\Delta}{2}])}{\Delta^2} \quad \text{when independence}$$

$$= \lim_{\Delta \rightarrow 0} \underbrace{\frac{1}{\Delta} \mathbb{P}(X \in [x - \frac{\Delta}{2}, x + \frac{\Delta}{2}])}_{\downarrow P_X(x)} \underbrace{\frac{1}{\Delta} \mathbb{P}(Y \in [y - \frac{\Delta}{2}, y + \frac{\Delta}{2}])}_{\downarrow P_Y(y)}$$

$$E[X Y] = \iint_K xy P_X(x) P_Y(y) dx dy \quad \text{for independence.}$$

$$= \int x P_X(x) \int y P_Y(y) dx dy$$

$$= \int x P_X(x) E[Y]$$

$$= E[X] E[Y]$$

$$X \sim N(\mu_1, \sigma_1^2) \Rightarrow X+Y \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2) \quad \text{when independence.}$$

$$Y \sim N(\mu_2, \sigma_2^2)$$

b/c

$$E[X+Y] = E[X] + E[Y]$$

$$\begin{aligned} \text{Var}(X+Y) &= E[(X+Y - E[X+Y])^2] \\ &= E[(X - E[X]) + (Y - E[Y])]^2 \\ &= E[(X - E[X])^2 + (Y - E[Y])^2 + 2(X - E[X])(Y - E[Y])] \\ &= E[(X - E[X])^2] + E[(Y - E[Y])^2] + E[2(X - E[X])(Y - E[Y])]. \end{aligned}$$

$$= \text{Var}(X) + \text{Var}(Y) + 2 \underbrace{E[(X - E[X])(Y - E[Y])]}_{= 0}.$$

$$\begin{matrix} \downarrow \sigma^2 \\ \sigma^2 \end{matrix} \quad \begin{matrix} \downarrow \sigma^2 \\ \sigma^2 \end{matrix} \quad = E[X - E[X]] E[Y - E[Y]].$$

$$= E[X] - E[E[X]] \dots$$

$$= E[X] - E[X] \dots$$

$$= 0$$

$$\cdot X + X = 2X \sim N(2\mu, 4\sigma^2), \quad X + Y \sim N(2\mu, 2\sigma^2)$$

\uparrow
independent

due to
cancellation

$$E[(2x - E[2x])^2] = 4E[(x - E[x])^2]$$

$$= 4\sigma^2$$

- - - - - - - - - - - - - - - - - - -

Empirical Risk Minimization

Go back to differential privacy:

- loss(h, x_i)
- $R(H; D) = \frac{1}{n} \sum_{i=1}^n \text{loss}(h, x_i)$

Mechanism: w/ x probability faithfully
w/ $1-x$ probability unfaithfully
random

Can attacker deduce x , based on data?

Assume they know: $\Pr(\text{"Dem"}) = \frac{1}{4}$.

Observe $\{D, R, K, D, R, R\} = S$.

If we know $\Pr(X \in [d - \frac{\Delta}{2}, d + \frac{\Delta}{2}] \mid \text{Data} = S)$.

$$a^* = \arg \min_d \Pr(X \in [d - \frac{\Delta}{2}, d + \frac{\Delta}{2}] \mid \text{Data} = S). \quad \left. \right\} \text{Bayes'}$$

$$d^* = \arg \min_d P_X(d \mid \text{Data} = S)$$

$$= \frac{\Pr(\text{Data} = S \mid X \in \dots) \Pr(X \in \dots)}{\Pr(\text{Data} = S)}$$

$$\text{For small } \omega = \frac{\Pr(\text{Data} = S \mid X = d) P_X(d)}{\Pr(\text{Data} = S)}.$$

$$\propto \Pr(\text{Data} = S \mid X = d) P_X(d)$$

$$\begin{matrix} \uparrow & \uparrow \\ \text{likelihood} & \text{prior probability} \end{matrix}$$

Probability Theory for ERM

DP problem with x prob, reveal true preference
($1-x$) prob, uniform random.

$$\Pr(\text{"Dem"}) = \frac{1}{4}, \quad \Pr(\text{"Rep"}) = \frac{3}{4} \quad \text{Observe } \{D, R, K, D, R, R\} = S.$$

$$\mathbb{P}(x \in [d - \frac{\Delta}{2}, d + \frac{\Delta}{2}] | \text{Data} = s) = p_x(x | \text{Data} = s)$$

pick $x^* = \underset{d}{\operatorname{argmax}} \mathbb{P}(x \in [d - \frac{\Delta}{2}, d + \frac{\Delta}{2}] | \text{Data} = s)$. Posterior

Bayes' Thm: $\mathbb{P}(x \in [d - \frac{\Delta}{2}, d + \frac{\Delta}{2}] | \text{Data} = s)$.

$$= \mathbb{P}(\text{Data} = s | x \in [d - \frac{\Delta}{2}, d + \frac{\Delta}{2}]) \mathbb{P}(x \in [d - \frac{\Delta}{2}, d + \frac{\Delta}{2}]) / \mathbb{P}(\text{Data} = s)$$

$$\propto \mathbb{P}(\text{Data} = s | x \in [d - \frac{\Delta}{2}, d + \frac{\Delta}{2}]) \mathbb{P}(x \in [d - \frac{\Delta}{2}, d + \frac{\Delta}{2}])$$

small Δ

$$\propto \underbrace{\mathbb{P}(\text{Data} = s | x = d)}_{\substack{\uparrow \\ \text{Likelihood}}} \underbrace{p_x(d)}_{\substack{\uparrow \\ \text{prior probability / belief}}}$$

Maximum Likelihood Estimate

$$\mathbb{P}(\text{Data} = s | x = d) = \prod_{i=1}^6 \mathbb{P}(\text{Data}_{(i)} = s_{(i)} | x = d).$$

$$\text{Data}_{(1)} = s_{(1)} \wedge \text{Data}_{(2)} = s_{(2)} \dots$$

\uparrow independence

$$\begin{aligned} \mathbb{P}(\text{Data} = "D" | x = d) &= d \cdot \frac{1}{4} + (1-d) \cdot \frac{1}{2} \\ &= \frac{1}{2} - \frac{1}{4}d \end{aligned}$$

$$\mathbb{P}(\text{Data} = "K" | x = d) = 1 - \mathbb{P}(\text{Data} = "D" | x = d).$$

$$= \frac{1}{2} + \frac{1}{4}d$$

Empirical Risk

$$\downarrow = \left(\frac{1}{2} - \frac{1}{4}d \right)^2 \left(\frac{1}{2} + \frac{1}{4}d \right)^4$$

$$\text{Risk}(d | s) = -\ln \mathbb{P}(\text{Data} = s | x = d)$$

$$= -2 \ln \left(\frac{1}{2} - \frac{1}{4}d \right) - 4 \ln \left(\frac{1}{2} + \frac{1}{4}d \right).$$

$\mathbb{P}(\text{Data}_{(i+1)} | \text{Data}_{(i)}, x = d)$

$$\text{so } \frac{d}{d\alpha} \text{R}(d | s) = -2 \frac{-\frac{1}{4}}{\frac{1}{2} - \frac{1}{4}d} - 4 \frac{\frac{1}{4}}{\frac{1}{2} + \frac{1}{4}d} = 0$$

Not independent case

\Downarrow

$$\frac{1}{2} \cdot \frac{1}{\frac{1}{2} - \frac{1}{4}\alpha} = \frac{1}{\frac{1}{2} + \frac{1}{4}\alpha}$$

$$\frac{1}{2}(\frac{1}{2} + \frac{1}{4}\alpha) = \frac{1}{2} - \frac{1}{4}\alpha$$

$$\frac{1}{4} + \frac{1}{8}\alpha = \frac{1}{2} - \frac{1}{4}\alpha$$

$$\frac{3}{8}\alpha = \frac{1}{4}$$

$$\alpha = \frac{2}{3}$$

11

prior belief

$$p_x(\alpha) = \begin{cases} 1, & 0 \leq \alpha \leq 1 \\ 0, & \text{else} \end{cases}$$

Relating to groupwork

$$\theta - 1 \quad \theta \quad \theta + 1$$

$$\begin{array}{c} x \ x \ \checkmark \sim \checkmark \\ \downarrow \\ \theta^* \end{array}$$

$$x_i \sim \text{Unif}[\theta - 1, \theta + 1]$$

$P(\theta | D)$: Posterior

$$P(D|\theta) : \text{likelihood.} = \prod_{i=1}^n P(x_i | \theta)$$

$$p_{x_i}(n) = \begin{cases} \frac{1}{2}, & x_i \in [\theta - 1, \theta + 1] \\ 0, & \text{else} \end{cases}$$

$$\text{if } x_1, \dots, x_n \in [\theta^* - 1, \theta^* + 1]$$

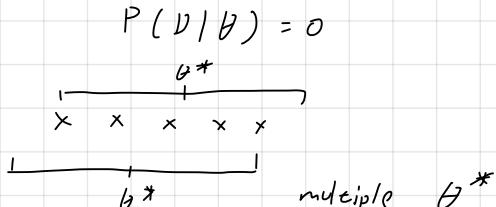
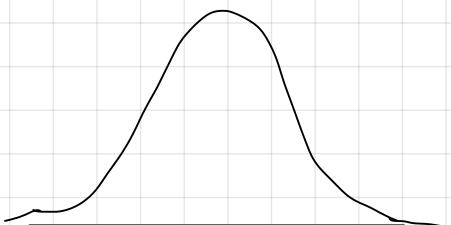
11

$$P(\theta | D) > 0$$

$$\text{if } \exists x_i \notin [\theta^* - 1, \theta^* + 1],$$

11

Central Limit Theorem :



$$P(D|\theta) = 0$$

$$\prod_{i=1}^n P(x_i | \theta) \propto \prod_{i=1}^n P(x_i | \theta) \cdot P(\theta) \rightarrow \text{Normal Distribution}$$

Bernoulli von Mises

Q: Why not Model data as normal R.V. ?

$$P(x_i | \theta) : x_i \sim N(\theta, \sigma^2).$$

$$\begin{aligned} \text{then } P(D|\theta) &= \prod_{i=1}^n P(x_i | \theta) \quad \text{likelihood} \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi} \cdot \sigma} e^{-\frac{(x_i - \theta)^2}{2\sigma^2}} \\ &\propto \prod_{i=1}^n e^{-\frac{(x_i - \theta)^2}{2\sigma^2}} \\ &= e^{-\sum_{i=1}^n \frac{(x_i - \theta)^2}{2\sigma^2}} \\ &= e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta)^2} \end{aligned}$$

Empirical risk

$$R(\theta, D) = -\ln P(D|\theta)$$

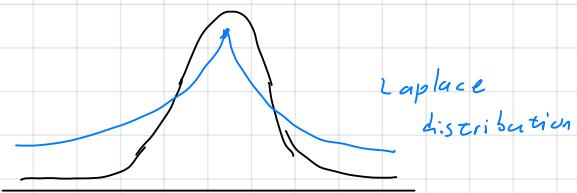
$$= \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta)^2$$

↗ Empirical risk squared loss

$$\theta^* = \arg \max_{\theta} P(D|\theta) = \arg \min_{\theta} R(\theta, D) = \frac{1}{n} \sum_{i=1}^n x_i$$

With finite data, how to model outliers?

Normal distribution



$$P(x_i | \theta) \propto \exp(-\frac{(x_i - \theta)^2}{2\sigma^2})$$

Value decreases faster

$$P(x_i | \theta) = \frac{1}{2} \exp(-|x_i - \theta|)$$

$$P(D|\theta) = \prod_{i=1}^n \exp(-|x_i - \theta|) \propto \exp\left(-\sum_{i=1}^n |x_i - \theta|\right)$$

$$R(\theta | D) = \sum_{i=1}^n |x_i - \theta|. \quad \text{Absolute loss}$$

Regression

$$y_i \approx H(x_i, \theta)$$

Linear regression

$$y_i \approx \langle \vec{x}_i, \vec{w} \rangle$$

$$y_i \sim N(\langle \vec{x}_i, \vec{w} \rangle, \sigma^2) \Rightarrow P(y_i | \vec{x}_i, \vec{w})$$

$$P(D | \vec{w}) = \prod_{i=1}^n P((x_i, y_i) | \vec{w})$$

$$\begin{aligned} P((x_i, y_i) | \vec{w}) &= P(y_i | \vec{x}_i, \vec{w}) \underbrace{P(\vec{x}_i | \vec{w})}_{\propto P(\vec{x}_i)} \\ &\propto P(y_i | \vec{x}_i, \vec{w}) \end{aligned}$$
$$\propto \exp(- (y_i - \langle \vec{x}_i, \vec{w} \rangle)^2 / 2\sigma^2).$$

Risk = Least square

P.3

A : J

B : J & k A

C : Heart

A : heart

B : face

C : red

Condition

$$\underline{P}(A \cap B | C) = \frac{1}{12}$$

$$\underline{P}(A \cap B | C) = \frac{4}{24} = \frac{1}{6}$$

$$\underline{P}(A | C) \underline{P}(B | C)$$

$$\underline{P}(A | C) \underline{P}(B | C) = \frac{12}{24} \cdot \frac{8}{24} = \frac{1}{2} \cdot \frac{1}{3} = \frac{1}{6}$$

$$= \frac{1}{12} \cdot \frac{4}{12} = \frac{4}{24} = \frac{1}{6}$$

$$\underline{P}(A | B \cap C) = \frac{4}{8} = \frac{1}{2}$$

$$\underline{P}(A | C) = \frac{1}{2}$$

$$\underline{P}(A | B \cap C) = \frac{1}{4}$$

$$\underline{P}(B | C)$$

$$\underline{P}(A | C) = \frac{1}{12}$$