Notation (I will try to be consistent!)

d: number of inputs w_i : weight from input *i* c: number of outputs w_{ii} : weight from *i* to *j* C_k : k^{th} output category *i*: input index N: number of patterns *j*: hidden unit index *n*: pattern label k: output unit index (not always!) t: target label a_i : weighted sum of inputs x: input data to unit *j*

g():activation function

y: output activation

z: hidden unit activation

· Unsupervised (A mobel of the data) • Supervised. (A mapping from input to targets) • Reinforcement. (learn from mistakes). · Initation

Regression:

There is a underlying line concepted by noise.

× 2 features

Linenr regression:

 $\begin{bmatrix} 1 & \chi_{i}^{\prime} & & & \\ 1 & & \chi_{i}^{N} \end{bmatrix} W \xrightarrow{\Lambda} targests$ $\begin{bmatrix} 1 & \chi_{i}^{\prime} & & & \\ 1 & & \chi_{i}^{N} \end{bmatrix} \xrightarrow{I}$ $\begin{bmatrix} 1 & \chi_{i}^{\prime} & & & \\ 1 & & \chi_{i}^{N} \end{bmatrix} \xrightarrow{I} \\ by \quad minimizes \quad the \quad Sum \quad Squareb \quad Error \quad (SSE).$

 $SSE = \frac{1}{2} \sum_{n=1}^{N} \left(\sum_{j=0}^{d} w_{j} x_{j}^{n} - t^{n} \right)^{2}$ $\int prediced$

number of examples/patterns. Solving for analytical solution: ... $\vec{w} = (x^T x)^{-1} x^T \vec{t}$. Gradient descent: .-- u^{weighted} sum of inputs $a^2 = w^7 \kappa^2$. y(x) = g(a)ケ Activation Partial derivative of error wrt weights $\frac{\partial \mu s E}{\partial w_{i}} = \dots = \frac{1}{N} \frac{N}{\sum} (y - z) \kappa_{i}$ SO_{j} $W_{i} = W_{i} + Q \frac{1}{N} \sum_{n=1}^{N} S^{n} \varkappa_{i}^{n}$ uphare rule. $= w_{i} + a \frac{1}{N} \sum_{n=1}^{N} (t^{n} - t^{n}) \lambda_{i}^{n}.$ Batch Learning Online Learning -> 5 GTD, charge weight immediately, for each data point. (train as hata comes in stream). у | Т.... | | | | Perceptions Trying to fit a function of Ð (a classifier). Binary threshold unit: December loyer $Y = \begin{cases} 1 & \text{if } \overline{Z} \ w; x; \ \overline{Z} \ \theta \\ \vdots & \ddots & \ddots \\ 0 & \text{otherwise.} \qquad Sum \ of \qquad weight. \end{cases}$ (X2) (X3) input layer. (\mathbf{x}) (regression). make a neurally - inspired machine that To could caregorize inputs and learn to bo this from example.

Perception Learning (perception convergence procedure).

=) w and vector (2 + - 2 ") is orchogonal. The distance l'is then: $\mathcal{L} = \frac{w^{T} \varkappa}{||w||} = \frac{-w_{o}}{||w||}, \text{ where } w = (w_{1}, \dots, w_{d}).$ l'erceptron as a linear discriminant. : · a two-class discriminat is one such that $y(x) = \alpha = w^T k$ χ in C₁ if $\alpha = 0$ no activation. else x in C2. · For multiple out put Y, Yc 0 0 0 unepare x is assigned Ck if k=argmax Y; (x) x, xb input U Xo n each region is convex sec => π^{4} , π^{13} , ... $= 7 n^{A}, n^{B} in R_{i}$ implies ant + (1-a) a B in K; Decision boundary is thus Yi(n) = Yi(n) for Ci and Ci, ... Logistic Kegression sigmoid function Ostimate of prob. of Cr or Cz. (still a Linear clossifier). $Y(x) = g(w^{7}x + w_{0})$ $g(x) = \frac{1}{1 + \rho^{-x}}$ logistic function as accivation. Derivation: if we assume our back is gaussian histribution (actually not in mond cases).

$$\begin{aligned} \mathbf{z}_{hl} \quad \mathbf{d}_{aca} \quad \mathbf{b}_{acchorism} \quad \mathbf{p}\left(n/c\right) &= \frac{1}{\sqrt{2\pi}} \frac{-\frac{1}{6\pi}}{6\pi} \left(\frac{\mathbf{z}_{h}}{\mathbf{z}_{h}}\right)^{2\pi}} \\ \quad \mathbf{p}\left(\mathbf{z}_{h}\right) \mathbf{z}_{h} &= \frac{1}{\sqrt{2\pi}} \frac{-\frac{1}{6\pi}}{6\pi} \left(\frac{\mathbf{z}_{h}}{\mathbf{z}_{h}}\right)^{2\pi}} \\ \quad \mathbf{p}\left(\mathbf{z}_{h}\right) \mathbf{z}_{h} &= \frac{1}{\sqrt{2\pi}} \frac{-\frac{1}{6\pi}}{6\pi} \left(\frac{\mathbf{z}_{h}}{\mathbf{z}_{h}}\right)^{2\pi}} \\ \quad \mathbf{k}_{g} \quad \mathbf{b}_{gres} \quad rale \quad \mathbf{p}\left(c, |\mathbf{x}\rangle\right) &= \frac{\mathbf{P}\left(\mathbf{z}_{h}(z_{h}) \mathbf{P}(c_{h})\right)}{\mathbf{P}(\mathbf{z}_{h})} \\ \quad \mathbf{p}\left(c_{h}|\mathbf{x}\rangle\right) &= \frac{\mathbf{P}\left(\mathbf{z}_{h}(z_{h}) \mathbf{P}(c_{h})\right)}{\mathbf{P}(\mathbf{z}_{h})} \\ \quad \mathbf{p}\left(c_{h}|\mathbf{z}\rangle\right) &= \frac{\mathbf{P}\left(\mathbf{z}_{h}(z_{h}) \mathbf{P}(c_{h})\right)}{\mathbf{P}(\mathbf{z}_{h}(z_{h}) \mathbf{P}(c_{h})} \\ \quad \mathbf{p}\left(c_{h}|\mathbf{z}\rangle\right) &= \frac{\mathbf{P}\left(c_{h}(z_{h}) \mathbf{P}(c_{h})\right)}{\mathbf{P}(\mathbf{z}_{h}(z_{h}) \mathbf{P}(c_{h})} \\ = \frac{1}{1 + e^{-\mathbf{R}}} \left(\frac{\mathbf{P}(\mathbf{z}_{h}|c_{h})\mathbf{P}(c_{h})}{\mathbf{P}(\mathbf{z}_{h}(z_{h}) \mathbf{P}(c_{h})}\right) \\ \quad \mathbf{z}_{h} &= \frac{1}{1 + e^{-\mathbf{R}}} \left(c_{h} \frac{\mathbf{P}(\mathbf{z}_{h}|c_{h})\mathbf{P}(c_{h})}{\mathbf{P}(\mathbf{z}_{h}(z_{h}) \mathbf{P}(c_{h})}\right) \\ \quad \mathbf{z}_{h} &= \frac{1}{1 + e^{-\mathbf{R}}} \left(c_{h} \frac{\mathbf{P}(\mathbf{z}_{h}|c_{h})\mathbf{P}(c_{h})}{\mathbf{P}(\mathbf{z}_{h}(z_{h}) \mathbf{P}(c_{h})}\right) \\ \quad \mathbf{z}_{h} &= \frac{1}{1 + e^{-\mathbf{R}}} \left(c_{h} \frac{\mathbf{P}(\mathbf{z}_{h}|c_{h})\mathbf{P}(c_{h})\right) \\ \quad \mathbf{z}_{h} &= \frac{1}{1 + e^{-\mathbf{R}}} \left(c_{h} \frac{\mathbf{P}(\mathbf{z}_{h}|c_{h})\mathbf{P}(c_{h})\right) \\ \quad \mathbf{z}_{h} &= \frac{1}{1 + e^{-\mathbf{R}}} \left(c_{h} \frac{\mathbf{P}(\mathbf{z}_{h}|c_{h})\mathbf{P}(c_{h})\right) \\ \quad \mathbf{z}_{h} &= \frac{1}{1 + e^{-\mathbf{R}}} \left(c_{h} \frac{\mathbf{P}(\mathbf{z}_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h})\mathbf{P}(c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|c_{h}|$$

$$= \frac{1}{\mu} \sum_{i=1}^{n} (i \cdot 3) (-3(a) \frac{da}{d^{2n}})$$

$$= \frac{1}{\mu} \sum_{i=1}^{n} (i \cdot 3) (-3(a)) x_{i}$$

$$= \frac{1}{\mu} \sum_{i=1}^{n} (b \cdot 3(a)) (x_{i})$$

$$W_{i} = W_{i} + \frac{d}{N} \sum_{i=1}^{n} (t^{n} - y^{n}) g(a) (1 - g(a)) x_{i}^{n}$$

$$= W_{i} + \frac{d}{N} \sum_{i=1}^{n} (t^{n} - y^{n}) g(a^{n}) (1 - g(a^{n})) x_{i}^{n}$$

$$g(a) = g(a) (1 - g(a))$$

$$W_{ise} \quad (reg - entropy - as - b) extrine - function.$$

$$Generalization \quad if \quad (agive - repressive to - multiple - output - g(a))$$

$$W_{ise} \quad (reg - entropy - as - b) extrine - function.$$

$$Generalization \quad if \quad (agive - repressive to - multiple - output - g(a))$$

$$W_{ise} \quad (reg - entropy - as - b) extrine - function.$$

$$Generalization \quad if \quad (agive - repressive to - multiple - output - g(a))$$

$$W_{ise} \quad (reg - entropy - as - b) extrine - output - g(a)$$

$$W_{ise} \quad (reg - entropy - as - b) extrine - output - g(a)$$

$$W_{ise} \quad (reg - entropy - as - b) extrine - output - g(a)$$

$$W_{ise} \quad (reg - entropy - as - b) extrine - output - g(a)$$

$$W_{ise} \quad (reg - entropy - as - b) extrine - output - g(a)$$

$$W_{ise} \quad (reg - entropy - as - b) extrine - output - g(a)$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entropy - g(a)) = \frac{e^{i}}{2}$$

$$W_{ise} \quad (reg - entr$$

. . .

· Create new feature from old feature for learning (old fashion). Formard Propagazion & Backwark Propogazion: Error (i) (i)Chain rule of grahest $(z_i = z_i) \xrightarrow{w_i} a_j$ Back prop. Learning: $W_{ij} = W_{ij} - \frac{\partial J}{\partial W_{ij}}$ (Zi mean Dutput of hibban unit or input). J is objective. is objective. $\frac{\partial J}{\partial w_{ij}} = \frac{\partial J}{\partial a_j} \frac{\partial a_j}{\partial w_{ij}} = \frac{\partial \overline{Z}_{l=0}^{k} w_{lj} \overline{z}_{l}}{\partial w_{ij}} = \frac{\partial \overline{Z}_{l}}{\partial w_{ij}} = \frac{\partial$ $= \frac{\partial \mathcal{J}}{\partial a_j} \mathcal{Z};$ 50 Z;=g[ai) → aj • Define $\xi_j = -\frac{\partial J}{\partial a_j} \bigg(= -\frac{\partial \chi_{pre}}{\partial a} = -\frac{\partial \chi_{pre}}{\partial y} \frac{\partial y}{\partial a} = (t-y) \bigg)$ logistic regression $\left(w_{ij} \leftarrow w_{ij} - \frac{\partial 7}{\partial w_{ij}} = w_{ij} - \frac{\partial 7}{\partial a_j} \frac{\partial a_j}{\partial w_{ij}} \mathbf{x}_i = w_{ij} + (t_j - f_j) \mathbf{x}_i \right)$ so we have $\frac{\partial J}{\partial w_i} = \frac{\partial J}{\partial a_i} \frac{\partial a_j}{\partial w_i} = - S_j Z_j$ By delta rule, wij \leftarrow wij $-\frac{\partial J}{J^{wij}} = w_{ij} + \delta_j Z_i$ (weight changes in proportion to the input and the belta at that connection). For output unit, delta rule is Wij < wij + (tj - 4;) Zi hidden unit, $\frac{\partial J}{\partial a_j} = \sum_k \frac{\partial J}{\partial a_k} \frac{\partial a_k}{\partial a_j}$ Ĩ-01 a_k a_k \dots a_k $\int \sqrt{w_{jk}}$ iben ; $Z_j = g(\alpha_j)$

$$\begin{aligned} s_{\mathbf{r}} \quad berivation: \quad \frac{\partial T}{\partial \sigma_{\mathbf{j}}} = \sum_{\mathbf{k}} \frac{\partial T}{\partial \sigma_{\mathbf{k}}} \quad \frac{\partial \sigma_{\mathbf{k}}}{\partial \sigma_{\mathbf{j}}} \\ &= -\sum_{\mathbf{k}} S_{\mathbf{k}} \quad \frac{\partial \sigma_{\mathbf{k}}}{\partial \sigma_{\mathbf{j}}} \quad \frac{\partial Z_{\mathbf{j}}}{\partial \sigma_{\mathbf{j}}} \\ &= -\sum_{\mathbf{k}} S_{\mathbf{k}} \quad \frac{\partial \sigma_{\mathbf{k}}}{\partial \sigma_{\mathbf{j}}} \quad \frac{\partial Z_{\mathbf{j}}}{\partial \sigma_{\mathbf{j}}} \\ &= -\sum_{\mathbf{k}} S_{\mathbf{k}} \quad \frac{\partial \sigma_{\mathbf{k}}}{\partial \sigma_{\mathbf{j}}} \quad \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{j}}} \\ &= -\frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{j}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad \sum_{\mathbf{k}} \frac{\partial w_{\mathbf{k}} 2}{\partial \sigma_{\mathbf{j}}} \quad every \quad even \quad is \quad O \\ &= -\frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{j}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad \sum_{\mathbf{k}} \frac{\partial w_{\mathbf{k}} 2}{\partial \sigma_{\mathbf{k}}} \quad every \quad even \quad is \quad O \\ &= -\frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{j}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{j}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{j}}} = \frac{\partial (\sigma_{\mathbf{k}})}{\partial \sigma_{\mathbf{j}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{j}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{j}}} = \frac{\partial (\sigma_{\mathbf{k}})}{\partial \sigma_{\mathbf{j}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{j}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{j}}} = \frac{\partial (\sigma_{\mathbf{k}})}{\partial \sigma_{\mathbf{j}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{j}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{j}}} = \frac{\partial (\sigma_{\mathbf{k}})}{\partial \sigma_{\mathbf{j}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{j}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{k}}} = \frac{\partial (\sigma_{\mathbf{k}})}{\partial \sigma_{\mathbf{j}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{j}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{k}}} = \frac{\partial (\sigma_{\mathbf{k}})}{\partial \sigma_{\mathbf{k}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{k}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{k}}} = \frac{\partial (\sigma_{\mathbf{k}})}{\partial \sigma_{\mathbf{k}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{k}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{k}}} = \frac{\partial (\sigma_{\mathbf{k}})}{\partial \sigma_{\mathbf{k}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{k}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{k}}} = \frac{\partial (\sigma_{\mathbf{k}})}{\partial \sigma_{\mathbf{k}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{k}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{k}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{k}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{k}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{k}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}}{\partial \sigma_{\mathbf{k}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{k}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{k}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{k}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}}{\partial \sigma_{\mathbf{k}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{k}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}}{\partial \sigma_{\mathbf{k}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{k}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}}{\partial \sigma_{\mathbf{k}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{k}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{k}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{k}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial \sigma_{\mathbf{k}}} \sum_{\mathbf{k}} S_{\mathbf{k}} \quad W_{\mathbf{k}} \mathbf{k} \\ \frac{\partial Z_{\mathbf{k}}}{\partial$$

(unline) Back prop : · present pattorn to network, feed forward until output. . compare & (usually t-J) · back prop & through network. Every unit has a delta. · update wij $w_{ij} \leftarrow w_{ij} + d\delta_j z_j$ Depende on initial condition, backprop gives many solution. (highly non-convex neework). Learn internal representation through hidden unles. symmetry Growth center ел) хок hikken _ υ $\hat{\mathbf{0}}$ 50 OUTPUT 0 0 0 ן ו ק ק hibben $\frac{-9}{\sqrt{+3}}
 \frac{-9}{\sqrt{+7}}
 \frac{(+7)}{\sqrt{-3}}
 \frac{(+7)}{\sqrt{-9}}$ ен) -3,7 (i) +3 PD) Net Talk Architecture , input uses a one-hot · averaging vector in tree structure (hibben unit).

Pm) Hinton's Family tree

Objective Function & MLE. likelihoud . privs Maximum Likelihoob. Want W that maximize $P(w|D) = \frac{P(0|w)P(w)}{P(D)}$ no-malizing constant IJ max p (w 1D) = max p (D/W) (maximize likelihosh How we noted the back $\mathcal{L} = \frac{N}{\Pi} P(x^{n})$ Find parameter by argmax L=> InL=In TT P(x^) Kob (stanbard distribution example) argmin - In TT p (x) = . _ . = ... Mobeling input-output data: = leads to setting IM to empirical menn. $\mathcal{L} = \frac{N}{\Pi} p(\mathbf{x}^{n}, t^{n}) = \frac{N}{\Pi} P(t^{n} | \mathbf{x}^{n}) P(\mathbf{x}^{n})$ $-\ln L = -\sum_{n=1}^{N} (\ln p(t^{n} | x^{n}) + \ln p(x^{n})).$ • Assume gaussian noise $h(x) \leftarrow model h. \qquad p(t^{n}/x^{n}) = \frac{1}{\sqrt{2\pi}6^{2}} e^{-\frac{(t^{n} - h(x^{n}))^{2}}{26^{2}}}$ $= \frac{1}{\sqrt{2\pi}6^{2}} e^{-\frac{(t^{n} - y(x^{n};w))^{2}}{26^{2}}} (for all both and a both a both$ The error (negative log likelihook) $\vec{E} = -\ln \frac{1}{\sqrt{2 - 1} - 6^{2} \ln \frac{N}{p_{2}}} - \frac{(e^{n} - y(\mu^{n}; m))^{2}}{2 - 6^{2}}$ (ross - Entropy: $Y(x^{n}) = P(c_{1}|z^{n})$ $= \frac{1}{26^{2}} \sum_{n=1}^{N} (t^{n} - y(x^{n}, w))^{2} + \ln \left(\int 2\pi 6^{-1} N \right).$ Constant Modeling two Categories is modeling coin flip likc sum squaret loss! Ų Bernoulli discribución V

$$p(e^{n}|x^{2}) = (5^{n})^{n} (1 + y^{2})^{(1 + q^{2})} = y^{n} (1 + y^{2})^{(1 + q^{2})} = y^{n}.$$

if $e^{n} = 1$, $p(e^{n}|x^{n}) = y^{n} (1 + y^{2})^{(1 + q^{2})} = \frac{y^{n}}{2n!} (5^{n})^{(1 + q^{2})}.$

 5_{0} , d_{0} likelikask; $L = \prod_{n=1}^{N} p(e^{n}|x^{n}) = \prod_{n=1}^{M} (5^{n})^{(1 + q^{2})}.$

 $-1nL = -ln \prod_{n=1}^{M} (1 + y^{2})^{(n-1)}.$

 $-1nL = -ln \prod_{n=1}^{M} (1 + y^{2})^{(1 + q^{2})}.$

 $-1nL = -ln \prod_{n=1}^{M} (1 + y^{2})^{(1 + q^{2})}.$

 $p(e^{n}|x^{2}) = \sum_{n=1}^{M} ln (s^{2})^{n} (1 + y^{2})^{(1 + q^{2})}.$

Entropy: n reduces of information in a message above a median variable.

 $-P(n) \log P(n)$

Entropy: n reduces of $P(x^{n}) \ln P(x^{n}).$

 $(ross - colored - f \prod_{n=1}^{M} P(e^{n}) \ln P(x^{n}).$

 $(ross - colored - 2 \prod_{n=1}^{M} P(e^{n}) \ln P(x^{n}).$

 $(ross - colored - 2 \prod_{n=1}^{M} P(e^{n}) \ln P(x^{n}).$

 $P(ch|x^{1}) = Y_{n}(x^{n}), \quad e_{h} = 1 + lf from rangeng, h, 0 order free

 $\Rightarrow p(e^{n}|x^{n}) = \prod_{n=1}^{m} (y_{h})^{(n)} = \dots = -\prod_{n=1}^{M} \sum_{n=1}^{m} e^{n} \ln x^{n}$

 $Samese Necont Normet.$

 y^{n}

 $f low from the colored or freedor appret.

if $x_{n} \neq x_{n}$ from same colored or freedor appret.

if $x_{n} \neq x_{n}$ from same colored or freedor appret.

if $x_{n} \neq x_{n}$ from same colored or freedor appret.

if $x_{n} \neq x_{n}$ from same colored or freedor appret.

The less function: $L(m_{n} = \frac{1}{2}, \frac{1}{2}, \frac{1}{2}) = (1 + \frac{1}{2}) =$$$

Generalization

Decl	W:n.	Over fiering	?						
. ,	Mark		12) and the					
	n i i i		۲ رو	uru x) take	difference	Crojos of	the data o.	, +esize	it.
• /	Regularization.	•			11	, ,			
	reture no	bel complexity	$(J = E + \lambda)$	c)					
	· L2 ->	minimize w 2	, weight sme	ller in	proportion	n to its	size.		
	· L,>	minimize W	, neight sma	ller at	a cons	stant rate.			
	. Minimize	$C = w ^2 / (1)$	w `+1), pena	lize bi	g neight	less while	e penalizing	small	neights more.
· /	Dropout								
=>	p robabilistica	ully turn off	some frac	tion ^D	f hidden	(mits,			
• Ē,	al statente								
	15 3 copping								
1 A	bb noise	to inputs/ne	ight / hibben	Unit act	vations				
	make mobel	more robust	to perturbatio	ns.					
Neura)	Mee Trick ;								
	5GD.			Batch	Learn	112			
	,								
	I. get	une example		1.	get one	example		mini-batch -7 by -u	leurning.
	2. 00 mpc	ite gradient		2.	compute	gradient		small bg	tch of data
	3. Change	e neight		3.	abh to	runing	avrage e	rach time	changing the nois be
	4. Zf se	ien all example,	roshuffle	4.	If seen	all examj	ole, chanse	neighe	
								effici	PAE
Shuf	fling diat	a before	partition t	to n	ini-batch	,		chue to	marmul.
				should	have	hiverse			
				VIG . 9	- grismiael	· · · · ·			
РCI	A								
i dea:	ex) ;f all	positive inpo	ie, wij = h	vij + do	fjxi				
				Same sig	n for a	Il weights.			
					· ·				

if Xi's all positive, then the weight changes are all tor are highly correlated variables en) bytimization =) get two hiffpene information. en) difference scale of inpure PCA -> shifes mean of input variable to 0. · accorrelatos the input. · throw away dimension with smaller elsowalues. (himensionalizy reduction). · divide by standark deviation, make them all of equal size. Mean PCA Z-scoring vs. Cancellation • Z-scoring shifts the mean of the input variables to be 0 - not all positive or Covariance negative! Equalizati • Z-scoring does *not* decorrelate the inputs - because it is applied to each variable independently. • And...you *cannot* throw away the dimensions with the smallest eigenvalues -because you don't compute the eigenvectors. • Z-scoring divides by the standard deviation, making all of the variables 0 mean and unit standard deviation. So the variables are all of the same size ٠ 1/18/24 Neural Networks/Deep Learning Different activation function

recommentale signoit: f(x) = 1.7159 * tanh(0.667 x).

problem. stanbark signoit makes all input to next lagar positive.

 $f(t/-1) = t/-1 \quad (if input has unit variance, on tput has unit variance).$

Weight Initialization:
Wright connect inscinize to 0
if initialize to 0 => 8 is of the hidden will all be the same.
Wright initialization contributed with:
Imput normalization
Choice of sympole.
Suppore X: norm 0, une variance. Wij man 0, on average, STO of aj = XW
mill be:
Var (XW) = Var (X) Var (W) + Var (X) (E(W)) + Var (W) (E(X)) +
= Var (W).
So, STO of aj =
$$\int Var (V) = \int \overline{Z}$$
, (Uj) +
To bee 6(aj) = 1, we see Wis to be $(\frac{1}{2}m)^{V2}$.
To bee 6(aj) = 1, we see Wis to be $(\frac{1}{2}m)^{V2}$.
Intervalues, if highter unit has big famin, small charger on many of its incoming
meights con cause learning weight when famin is big, and vice-versa.

Bacch Normalization

iden: normalizes all inputs and inpurs to each hidden layer throughout the network.

(on a per unit basis, over each min; bacch).

1. Z-score each inpue variable over mini-bach.

 $\hat{a}_{i} = \frac{a_{i} - u_{i}}{(6_{i}^{2} + 6)}$ where u_{i} , b_{i}^{2} are mean and variance of a_{i} .

ner inpur to a unit are any layor of the normark.

2. Allow nearth laws a who if necessary

$$\widehat{\alpha_{i}} = Y, \widehat{\alpha}_{i} + \delta_{i}$$

$$T$$

$$Iconsolv property (ai = bar, a = a_{1...m});$$
Parameters to be lamin-batch: $B = \{a_{1...m}\};$
Purput: $\{\widehat{\alpha_{i}} = BN_{Y,\beta}(a_{i})\}$

$$\mu_{B} \leftarrow \frac{1}{m} \sum_{i=1}^{m} \alpha_{i} \qquad // \text{mini-batch mean}$$

$$\sigma_{B}^{2} \leftarrow \frac{1}{m} \sum_{i=1}^{m} (a_{i} - \mu_{B})^{2} \qquad // \text{mini-batch variance}$$

$$\widehat{\alpha_{i}} \leftarrow \frac{a_{i} - \mu_{B}}{\sqrt{\sigma_{B}^{2} + \epsilon}} \qquad // \text{normalize}$$

$$\widehat{\alpha_{i}} \leftarrow \frac{a_{i} - \mu_{B}}{\sqrt{\sigma_{B}^{2} + \epsilon}} \qquad // \text{normalize}$$

$$\widehat{\alpha_{i}} \leftarrow \frac{a_{i} - \mu_{B}}{\sqrt{\sigma_{B}^{2} + \epsilon}} \qquad // \text{normalize}$$

$$\widehat{\alpha_{i}} \leftarrow \widehat{\gamma_{a}} + \beta = BN_{\gamma,\beta}(a_{i}) \qquad // \text{scale and shift}$$

$$M \text{onersum}$$

$$ueigle charge are size c.$$

$$D \le V \ge V \ge \sqrt{c_{-1}} - \alpha = \frac{2\epsilon}{2m}(c)$$

$$province and (constant) = \frac{2\epsilon}{2m}(c)$$

$$province and (constant) = \frac{2\epsilon}{2m}(c)$$

$$province and (constant) = \frac{2\epsilon}{2m}(c)$$

$$how is behave:$$

$$D = (a + constant) = \frac{2\epsilon}{2m}(c)$$

$$how is behave:$$

$$D = (a + constant) = \frac{2\epsilon}{2m}(c)$$

$$how is behave:$$

$$D = (a + \frac{2\epsilon}{2m}) \quad (by scherz) = constant).$$
Necessor Mechals:

$$fine mode a = big jump in director of previous Accumulated gradient.$$

$$dis magnet are gradient when yow end male a correction.$$

 $\begin{array}{l} g_{ij}(t) = g_{ij}\left(t-1\right) + 0.05. \\ \left(\begin{array}{c} \frac{\partial E}{\partial w_{ij}}\left(t\right) \frac{\partial E}{\partial w_{ij}}\left(t-1\right) \right) & 70 \end{array} \\ \end{array}$ $\begin{array}{l} & \uparrow \\ & \downarrow \\ & \downarrow \\ & \uparrow \\ & \downarrow \\ \\ & \downarrow \\ \\ & \downarrow$ Apaptice learning race: $Dwij = - \xi g_{ij} \frac{\lambda E}{\lambda^{wij}}$ - decrease otherwise. Slobal local learning arkjusement. rate $g_{ij}(\tau) = g_{ij}(\tau - 1) \cdot 0.95.$ =) big gain being rupibly when oscilliation starts. Apapeire learning race only heals with axis-alignet effects. Apapeire learning rate affect individual weights , which charge the learning rate along that axis. Resilient Backprop (Rprop): Only for full-batch learning: · increase step size for a weight multiplicatively (times 1.2) if the signs of its last two gradients agree · otherwise, decrease the step size multiplicatively (times 0.5) Does not work for mini-back learning: Kprop is to use the gradient but also divide by the size of the gradient $\left(sign(x) = \frac{x}{|x|}\right)$, problem is that we divide by a different number for each mini-bach. Solution - RMSprop: Keep a moving average of squarely gradient for each weight. Mean Square (w, t) = 0.9 Mean Square (w, t-1) + 0.1 (f E/fw(c))²

Then, divide the gradient by Menn Square (U, t).

Convolutional Network (CNN). Problem : · Large images size (1140 x 648) Propercies of f. Nearby pixel correlaze with nearby pixel: locality Rent Worth / Image. Statistics of pixel are relatively uniform: stationary statistics · object doesn't depend on its location in image.: translation invariance , object make of parts. Compositionality Why CNN useful to solve above publics? \mathbb{U} . "Locality": small, local receptive field and learneh features (kernel). · "Scationary stacistics" : replicated across the image (if a feature is useful in one place, it's useful in others). · "Translation invariance": syntial pooling. · Objects are made of parts: receptive fields goe larger deaper in the nec. size of visual world that innervaces the neuron field Visual feature in visual world that most stimulate neuron.



Activation Maps pooling ---> · · · => · 2×2 set of response to I number (Translation Invariance). (1/4 size) Convolución lenús to shift equivariance. (shifting number in difference location does not affect the feature map). Original image Feature map for 1 feature 2 If we move the 3 over... the features just shift over \mathbb{Z} This is shift *equivariance* Overview of Conunets. I. Image Nee Data · Fred-forward 1 Alex Nee 8 layers - Convolve inpuz. VGG 19 layers \mathcal{J} - Non-linearitz (rectified linear) - Pooling (local max). U Google Nee 22 layers · Sceperviset 1 Ţ · Convolution filter trained by back prop. 1 ResNet 152 layers ١ Images -> Convolution -> Non-linear -> Pooling -> Feature Maps (filter)

Nores From Stanford (5231n. CNN (softmax). ex) (12 filters) (software) INPUT -> CONV -> RELU -> POOL -> FC [32×32×3] [32×32×12] [32×32×12] [16×16×12] [1×1×10] RELU POOL FC CONV cont Gin Х 12 Granetes X \times hyperparameter \succ \succ RELU RELU RELU RELU RELU CONV CONV CONV CONV CONV ↓ ↓ ↓ ↓ ↓ ↓ airplane ship horse Convolutional Layer · cach filter will produce a 2-D activation map. · Stacking activation map along the depth dimension probleme the output volume. · connece each neuron to only a local region of the input volume (receptive field). · the extend of connectivity along the depth axis is always equal to the kepth of input volume. en inpue volume [32×32×3] receptive field (filter size) [5x5x3] then a total of 5×5×3 = 75 weights.

en) input volume [16 × 16 × 20].

receptive fielde [3×3×20].

then a total of $3 \times 3 \times 20 = 180$ weights.

Three hyperparameters to control output volume $\begin{array}{c|c} x_0 & u_0 \\ \hline \\ axon from a neuron \\ w_1x_1 \\ \hline \\ w_2x_2 \\ \end{array} \begin{array}{c} \text{cell body} \\ f \\ \text{cell body} \\ f \\ \text{output axon} \\ \text{function} \\ \hline \\ \\ \text{function} \\ \end{array}$ | · Depth: # of filters (K) Same neight. 1 · Stribe: # of jumps when filter noves · Zero-pahainz: whether to pad input · W: input volume size volume with O around the border. · F: receptive field size - - - - -· S; seride 3×3 filter · P: amour of zero pathing $= \frac{7 - 3 + 2 \cdot 0}{1} + 1$ Then surprise volume size is = 4+1 = 5 $\frac{\left(W-F+2P\right)}{5}+1$ =7 5 × 5 OUEPUE. 7×7 choice of Stribe is important stribe = 1, pah = 0 to ensure that the above division _ _ _ _ _ _ _ _ _ Yielks an integer. Paranter sharing. Thus, to have input and output the same size, we set zero padiding since we are using filter striking over the inages, the # of weights will be p = (r - 1)/2F × (# of filter) × (depth of output). this is quarantee by translationally invariance, if detecting a feature is important at some location, it should be useful at some other location as well.

If learning difference feature on difference location of images is necessary, den its common to relax parameter sharing schema J. Locally - Connected Layer. Lonv a complete visualization Input Volume (+pad 1) (7x7x3) Filter W0 (3x3x3) Filter W1 (3x3x3) Output Volume (3x3x2) x[:,:,0] 0 0 0 0 0 0 0 0 w0[:,:,0] w1[:,:,0] 1 -1 0 o[:,:,0] 4 3 2 0 1 0 2 1 2 0 -1 0 -1 0 -1 0 2 1 7
 0
 0
 1
 0
 1
 0
 0

 0
 1
 1
 0
 0
 0
 0
 0 1 0 1 -1 0 -2 0 3 w1[:,:,1] 0 1 0 o[:,:,1] 1 4 0 ₩θ[:,:,1] 0 -1 Ø 0 0 0 2 2 2 0 -100 -2 7 6 1 1 -1 0 0 0 0 0 0 0 1 1 1 1 1 1 2 0 3 0 0 0 0 0 0 0 w1[:,:,2] 1 1 -1 w0[:,;,2] x[:,:,1] 0 0 0 0 0 0 0 -101 0 1/-1 -1 0 0 0 1 2 1 1 0 0 0 2 1 2 0 1 0 0/1 -1 0 -1 0
 0
 1
 2
 1
 1
 0

 0
 1
 0
 2
 2
 1
 0

 0
 1
 1
 2
 1
 1
 0
 Bias b0 (1x1x1) Bias b1 (1x1x1) b0[*,:,0] b1[:,:,0] 0 0 0 0 0 0 0 x[:,:,2] 0 0 0 0 0 0 0 0 2 0 2 0 1 0 0 0 2 0 1 0 0 0 2 0 1 0 0 0 2 0 0 1 0 0 0 2 0 0 1 0 toggle movement 9 1 0 9 0 0 0 0 0 2 0 0 0 0 0 0 0 0 0 0 0 Pooling layer · progressively reduce the spatial size of the representation to reduce the parameters. · W, H, D, representing input volume size. · For pooling layer, it requires two hyperparameters: · Spatial Extent F · stille S Then output volume: $W_2 = (W_1 - F)/5 + 1$ H2 = (H, - F)/S +1



 $D_2 = \mathcal{D}_1$

Pooling layer downsamples the volume spatially, independently in each depth slice of the input volume. Left: In this example, the input volume of size [224x224x64] is pooled with filter size 2, stride 2 into output volume of size [112x112x64]. Notice that the volume depth is preserved. Right: The most common downsampling operation is max, giving rise to max pooling, here shown with a stride of 2. That is, each max is taken over 4 numbers (little 2x2 square).

Convnet Architectures

Most common form:

INPUI -> [[[UNV -> RELU] × N -> POOL ?]×M -> [FC -> RELU]×K -> FC



Filtering · [Wright learned hurning training) · each unit looks at local winkow. · same feature computely every-here. Mon-linearien . · ReLU (applieh per-pixel). · Output = max (0, input). Spacial Pooling . max is besc · sum or avorage uset less · applied to each filte layer separately. · Scribe: how for to slike over each image role: invariance to small transformation. larger receptive field. Component of each laya: pixol/feature => filter with learnoh dictionary => Non-linearity \mathcal{P} Spacial local max pooling ψ · Barch normalization. Output features. How to choose architecture ? · Hyperparameter · en) # layers, # fecture majos, parans of fecture maps.



Peep Model: iben: design one modular stracture, then juse all more of them 3×3 conv lager better than single 7×7 b/c non-linearlists. · Network - in - network : add two per-pixel fully connected lagers before pooling. (onv -> 1×1 (onv -> 1×1 (onv N;N 13/ock · global average pooling. replace final fully connected layer: one feature map por category. · What is IXI COAN? [36×36×64] -> [36×56×32]. 1×1 Lonv with 32 filters (each filter size (1×1×64), performs a 64-dim dot product). => essentially a dot produce of filter wishes and 64 chandles abrantage: complete control of computation !!! Residual Neework idec: introduce "pass through" into each layer. unly residual needs to be learned. (the nervoile learn X idencity. how much to change X) $F(n) = H(n) - \kappa$ × -> neight layer -> neight layer -> (+) -> relu. F(x) F(x) + x = H(n)

Visualization for higher layer nerwork:



· max activation from feature map associated with each filter. · deconvaec to project back to pixel space. · pooling switch.

Receptive field of the first layer is the filter size
 Receptive field (w.r.t. input image) of a deeper layer depends on all previous layers' filter size and strides

Correspondence between a feature map pixel and an image pixel is not unique
Map a feature map pixel to the center of the

receptive field on the image in the SPP-net paper

Training Tips · Annealing learning rate. (> Starz lasse, slouly reduce. . Visualize Spature map.

GOOD BAD BAD too noisy too lack structure correlated

Good training: learned filters exhibit structure and are uncorrelated.

· Training liverges -> decrease learning race · Network underperforming -> make nee larger.

Recurrente Neewoode

Represent time in network ? map time into space map time into state of the network. Time -> Space. We-2 inpuz (z-2) inpuz (z-1) inpuz (z) · Auto regressive: · Feed-forward net: input(t-2) input(t-1) h.dben input (t) Tine -> State · use activation memory: new things are processed basely upon what has come before. (network state). input -> hibben -> output => O Jorhan net input ____ hibben ____ output & Elman new (simple recurrent net). onepne -> hibben -> outpore

equivolently, Lone-stop backprop through zime). input ->> hibben ->> output 1-1 way back hicken(t-1) (previous searce). Unrolling networks in time. - Dave Rumelhart · calculate total grabient over time time=3 (A) (B) (c)as an average and then modify **W**2 W3 WA **W**2 the weights. time=2 (A) **B** (c)W2 time=1 (A) В \bigcirc The recurrent net is just a W3 W **W**2 · propagate activation formark in time layered net that keeps reusing the same weights. time=0 (A) В \bigcirc (ZA(t), ZB(t)...) ZA(t.1) Deep Learning: RNNs 33 · propagace delta bademark. $(e_{\mathcal{X}}) \quad W_{3} = d \frac{1}{3} \sum_{k=1}^{3} \delta_{\mathcal{B}}(\epsilon) Z_{\mathcal{A}}(\epsilon-1)$ · upliate weight, Harli to train b/2 backwark pass is linear. problem of explohing / vanishing gradient: -> · switch logistic to ReLU. · weights big, gradient shrinks exponentially. ' In KIVIV, we can: · weights big, gradient grow exponentially. ' · explohing grudient. limic length of grahiene vector (gradient clipping) We can specify inputs in several ways: - Specify the initial states of all · Vanishing gradient: the units. w skip connection Specify the initial states of a subset of the units (the rest ↑ time are hiddens) w3 w4 **W**2 Specify the states of the W1 same subset of the units at every time step. I.e., we have inputs at every time step. W3 W4 **W**2 W1





Memorz is a matrix of linear neurons.

Transformer Neework Map time into space using shared weights communicate via attention $\langle | \rangle$ when you focus on particular aspect of your env. · Overt attention: more your eyes to work something · Covert attention: hirect your attention to something without moving your eyes How to direct attention in Neural Noc? J Sofemax + multiplicative connections. (Lontent - based Addressing). $W_{i}^{c}(i) \leftarrow \frac{e \times p(B \times [k_{\ell}, M_{\ell}(i)])}{\overline{Z_{j}} e \times p(B \times [k_{\ell}, M_{\ell}(j)])} \xrightarrow{e_{\pi}} M = \frac{3}{1} \frac{1}{2} \frac{4}{3} \frac{1}{4}$ $= \frac{3}{1} \frac{2}{2} \frac{3}{2} \frac{4}{2} \frac{1}{2} \frac{3}{2} \frac{4}{2} \frac{1}{2} \frac{3}{2} \frac{4}{2} \frac{1}{2} \frac{3}{2} \frac{4}{2} \frac{1}{2} \frac{1}{2} \frac{3}{2} \frac{4}{2} \frac{1}{2} \frac{1}{2} \frac{3}{2} \frac{4}{2} \frac{1}{2} \frac{1}{$ attention weight: w= (0.998, 0.001, 0.001). reak re < Z weli) Meli) 50, re=(2.997, 1.007, 3.997, 1.010) ~(3、1,4、1) · key vaces was computed by the network. (dynamic, network know what is was looking for). · Attention weight filterah the memory row via multiplication. · Match b/t key and memory deaner have to be perfect. Why transformer? _____> decoder needs different information at different time steps. · process each position in parallel. Attention-based methods -> problem with recurrence.

No parallelization since hidden states all depende on previous hidden states.





Generative Pre-Training 2.0 (GPT 2.0) (> trainch to predict next work in a massive dataset. Image Transformer · Takes in patches of image (no convolution) what it is learning? Vision Transformer (ViT) Class Bird Ball Car MLP Head encoher only + classifier. Attention Input Transformer Encoder Patch + Position Embedding -> 0: 0 2 3 4 5 6 7 8 8 * Extra learnable [class] embedding Linear Projection of Flattened Patche Universal Transformer Network: (7 shared weights between different layers (un rollede recurrent network - same computation is performeth over and over). Each comer can becilie abaptively how Universal Transformer Shared weights Networks: to run. many iteration between every tower beeneen layers. 6 times Shared weights between the different rows of acqua nets 2/18/25

RL, Peep Nee and Atori Agent learns ce police Tr Leurning to act in a world. $TT(5) = P(a_1, a_2, \cdots , a_n) >$ State st function from states to action probabilities. raise prob. of action a; when $(S_{t}) \xrightarrow{r_{t+1}} (S_{t+1}) \xrightarrow{s_{t+1}} (S_{t+1}) \xrightarrow{r_{t+2}} (S_{t+2}) \xrightarrow{r_{t+3}} (S_{t+3}) \xrightarrow{r_{t+3}} (S_{t+3})$ it was a gook move, lower ic if it was a back move " ex) PONG Policy network ram pixel Network (previous - current in age) 0 state(s) 0 action $policy(\pi)$ T(s) = P("up"|s)At the end, we use sign of the rewark on the gradience and apply it to every gradient we produced over the whole game. => win -> encourage actions that led to win lose -> discourage actions that led to loss. Policy gradient: 1) At each step of plan, sample from the softmax distribution at output: TT (S) = (0.1, 0.02, 0.6, ..., 0.01) network, mapping from states to probabilities of action

(2) Treat the sample as "teacher" $t = (0, 0, 1, 0, \dots, 0)$ for state / aution pair. 3 Computer weight change, adde then to a running average of weight change (some form of "gradient") (4) Multiply weight change by the sign of remark. 3 (hange the weights after one game. TD-Gammon and Alpha Go goal of learning: maximizing long-term expected rewark. Let re be remark at time t. $\binom{2}{maximize} \qquad R_t = r_{th} + \frac{1}{2}r_{trz} + \frac{1}{2}r_{zrz} + \cdots = \sum_{k=0}^{\infty} \frac{y^k}{r_{trk+1}}$ 0 ≤ Y ≤ 1 is the discount rate. $\langle \rangle$ ensures that expected reward converges. Policy: The (S, a) = prob. that ac = a when SE = S. State should satisfy Markov Property: Molel-baseh and Model-free.

Value Function:

 $V^{\pi}(5) = expected long-tarm return of being in this state, following policy Tr.$

TD - Grammon.

· first application of learning volue function using NN function approximator. representation of board as inpuc · learned value of board position as output. · neightet uplate using temporal difference rule on every move. currene estimate updatch to be closer to next hourd's value. expected value of boach position. (will be minimized when YEHI - it = 0).

Alpha - Go · Train two policy neeworks by supervised training on expere positions. one is shallow and fase one is dreep and accurate $DW = \frac{\partial t \ln(9)}{\partial w} \operatorname{sign}(r)$. train another network by policy gradient method. used for rollowes After superliset training, network is improved by playing isself.

A 3-h nerwork: value network trained to produce winner from every state.

from every state seen busing play, train to prehice win rate.

masking then Maskoh Autoencoluer reconstruct in par in pre-training. Description of the pre-training, use a simple decodier. BERT Step O: break image with non-overlapping parches, minimum decoker (2): mask 75 % of parches. (3): positional encoding of the patches. (4): learns to reconstruct the whole image. · self-supervised learning without any examples. · self-supervised learning with few data argumentation. AGI (advance ML/AI) · agent learn as much as they can about the world without interaction · differentiable architectures for planning and reasoning (by observation). · learn mulei-level absertion through long-term predicion and long-term planning.