

$f: \mathbb{R}^n \rightarrow \mathbb{R}$   $a \in \mathbb{R}^n$   
 $\frac{df}{dx}(a) \in \mathbb{R}^{1 \times n}$  (row vector)

**Taylor's approximation:**

$f(x) \approx f(a) + \left[ \frac{df}{dx} \right]_{x=a} (x-a)$   
**Gradient is derivative transposed**  
 $f(x) \approx f(a) + (\nabla_x f(a))^T (x-a)$

for matrices:  
 $f(x) \approx f(a) + \nabla_x f(a) \cdot (x-a)$   
 $f(x) \approx f(a) + \text{Tr}(\nabla_x f(a)^T (x-a))$   
 $[\nabla_x f(x)]_j = \frac{\partial f(x)}{\partial x_j}$   
 $\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} & \dots & \frac{\partial f(x)}{\partial x_n} \end{bmatrix}$

**Hessian: derivative of the gradient**

$[\nabla_x^2 f(x)]_{ij} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j} = \frac{\partial}{\partial x_j} [\nabla_x f(x)]_i$

**best quadratic approximation:**

$f(x) \approx f(a) + [\nabla_x f(a)]^T (x-a) + \frac{1}{2} (x-a)^T [\nabla_x^2 f(a)] (x-a)$

**derivative chain rule**

$\frac{d(f \circ g)(x)}{dx} = \left[ \frac{df}{dz} \right]_{z=g(x)} \cdot \left[ \frac{dg}{dx} \right]$

**Linear Regression**

examples: Gaussian Linear Regression  $p(y|x) = \mathcal{N}(y|w^T x, \sigma^2) \rightarrow$  OLS  
 $\frac{d \log p(y|x)}{d w} = \frac{1}{\sigma^2} (y - w^T x)$   
 $\frac{d \log p(y|x)}{d \sigma^2} = -\frac{1}{2\sigma^4} (y - w^T x)^2$   
 $\frac{d \log p(y|x)}{d \sigma} = -\frac{1}{\sigma^3} (y - w^T x)^2$   
 $\nabla_x^2 \log p(y|x) = -\frac{1}{\sigma^2} x x^T$   
 $\nabla_x^2 \text{LSE}(x) = -\frac{1}{\sigma^2} \text{diag}(x_1^2, \dots, x_n^2)$

**Miscellaneous**

$x^T A y = \sum_i \sum_j A_{ij} x_i y_j$   
 $\left[ \begin{matrix} x^T \\ A \\ y \end{matrix} \right]^T = x^T A^T y$   
 $L(\tilde{z}, \tilde{\lambda}, \tilde{\mu}) = f_0(\tilde{z}) + \sum \lambda_i f_i(\tilde{z}) + \sum \mu_j g_j(\tilde{z})$   
 $P^* = \min_{x \in \mathbb{R}^n} f_0(x)$   
s.t.  $f_i(x) \leq 0$   
 $g_j(x) \leq 0$

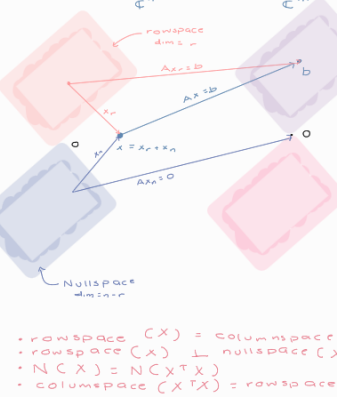
**Positive Semidefiniteness**

$A \in \mathbb{R}^{n \times n}$  symmetric  
 $\forall x \in \mathbb{R}^n \quad x^T A x \geq 0$   
 $\forall \lambda: \lambda \geq 0$  s.t.  $A = \sum \lambda_i v_i v_i^T$   
 $\exists U \in \mathbb{R}^{n \times n}$  s.t.  $A = U \Lambda U^T$

Properties of PSD matrices:  
 $A \succ 0, B \succ 0 \Rightarrow \alpha A + \beta B \succ 0$   
 $A \succ 0 \Rightarrow A^{-1} \succ 0$   
 $A \succ 0, B \succ 0 \Rightarrow \text{tr}(AB) \geq 0$   
 $A \succ 0, B \succ 0 \Rightarrow \text{tr}(AB) = 0 \iff AB = 0$   
 $\|A\|_F = \sqrt{\sum \lambda_i^2}$   
 $\lambda_{\max}(A) = \max_x \frac{x^T A x}{\|x\|_2}$

**Fundamental Theorem of Linear Algebra**

matrix  $A$  maps a vector from  $\mathbb{C}^n \rightarrow \mathbb{C}^m$



row space  $(X)$  = column space  $(X^T)$   
row space  $(X) \perp$  nullspace  $(X)$   
 $N(X) = N(X^T X)$   
column space  $(X^T X) =$  row space  $(X)$

**MLE**

goal: find a hypothesis model that maximizes the probability of the data  $\mathcal{D}$  parameterize the set of hypothesis models with  $\theta$   
 $\hat{\theta}_{MLE} = \arg \max_{\theta} \mathcal{L}(\theta; \mathcal{D}) = p(\text{data} = \mathcal{D} | \text{true model} = \theta)$   
 $= \arg \max_{\theta} \log \prod_{i=1}^n p(x_i; \theta)$   
 $= \arg \max_{\theta} \sum_{i=1}^n \log(p(x_i; \theta))$   
**Gaussian log-likelihood**  $y_i | x_i \sim \mathcal{N}(h_{\theta}(x_i), \sigma^2)$   
 $\hat{\theta}_{MLE} = \arg \max_{\theta} -\sum_{i=1}^n \frac{(y_i - h_{\theta}(x_i))^2}{2\sigma^2} - n \ln(\sigma)$   
 $= \arg \min_{\theta} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2$   
 $\rightarrow$  now, in the case of regression,  $h_{\theta}(x_i) = x_i^T \theta$   
 $\hat{\theta}_{MLE} = \arg \min_{\theta} \sum_{i=1}^n (y_i - x_i^T \theta)^2 = (X^T X)^{-1} X^T y$   
 $\rightarrow$  which is just OLS!

**MLE properties:**  
- consistent (more data  $\rightarrow$  convergence of true  $\theta$  value for  $\mathcal{D}$ )  
- statistically efficient (least variance parameter estimates)  
- value of  $p(\mathcal{D} | \theta_{MLE})$  is invariant to reparameterization  
eg  $\mathcal{N}(x|\mu, \sigma)$  vs  $\mathcal{N}(x|\mu, \sigma^2)$   
- can yield a parameter estimate even when data wasn't generated from family (which captures uncertainty)  
- only get a point estimate of parameter instead of distribution (MLE equivalent to maximizing cross-entropy/minimizing relative entropy)

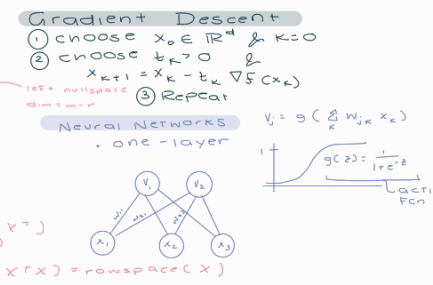
**Linear Regression**  
- consistent (more data  $\rightarrow$  convergence of true  $\theta$  value for  $\mathcal{D}$ )  
- statistically efficient (least variance parameter estimates)  
- value of  $p(\mathcal{D} | \theta_{MLE})$  is invariant to reparameterization  
eg  $\mathcal{N}(x|\mu, \sigma)$  vs  $\mathcal{N}(x|\mu, \sigma^2)$   
- can yield a parameter estimate even when data wasn't generated from family (which captures uncertainty)  
- only get a point estimate of parameter instead of distribution (MLE equivalent to maximizing cross-entropy/minimizing relative entropy)

**MAP for linear regression using Gaussian prior**  
- zero mean prior  $p(w) = \mathcal{N}(w; 0, \lambda I)$   
- Bayesian posterior  $p(w | \mathcal{D}) = p(\mathcal{D} | w) p(w)$   
 $w_{MAP} = \arg \max_w \log p(\mathcal{D} | w) p(w) = \arg \max_w \log p(\mathcal{D} | w) + \log \mathcal{N}(w; 0, \lambda I)$   
 $= \arg \max_w \sum_{i=1}^n \log \mathcal{N}(y_i; w^T x_i, \sigma^2) + \sum_{i=1}^d \log \left( \frac{1}{\sqrt{2\pi\lambda}} \exp\left(-\frac{w_i^2}{2\lambda}\right) \right)$   
 $= \arg \min_w \frac{1}{2\sigma^2} (y - Aw)^T (y - Aw) + \frac{\lambda}{2} w^T w$   
 $= \arg \min_w (y - Aw)^T (y - Aw) + \frac{\lambda}{\sigma^2} \|w\|_2^2$   
 $w_{RR} = (A^T A + \lambda I)^{-1} A^T y$   
 $\lambda > 0 \Rightarrow (A^T A + \lambda I)$  invertible

**Softmax Regression: classification**  
logistic regression  $\rightarrow 2$  classes  
 $p(y=1|x) = \frac{1}{1 + \exp(-z)}$   
 $p(y=0|x) = \frac{\exp(-z)}{1 + \exp(-z)}$   
**LOSS FNS:**  
(A) Squared error:  $L(z, y) = (z - y)^2$   
(B) Absolute error:  $L(z, y) = |z - y|$   
(C) Cross-entropy:  $L(z, y) = -y \ln(z) - (1-y) \ln(1-z)$   
**COST FNS to minimize:**  
(1) Mean loss:  $J(w) = \frac{1}{n} \sum L(w(x_i), y_i)$   
(2) Max loss:  $J(w) = \max L(w(x_i), y_i)$   
(3) Weighted sum  $J(w) = \sum w_i L(w(x_i), y_i)$   
(4) L2 penalty:  $J(w) = \sum w_i^2 + \lambda \|w\|_2^2$   
(5) L1 penalty:  $J(w) = \sum |w_i| + \lambda \|w\|_1$

**Famous Regression Methods**  
- Least-sq. linear regr:  $(1) + (A) + (4)$   
- weighted Least-sq. linear regr:  $(1) + (A) + (3)$   
- ridge regression:  $(1) + (A) + (4)$   
- LASSO:  $(1) + (A) + (5)$   
- Logistic regression:  $(2) + (C) + (5)$   
**Quadratic Program**  
- complex cost  $\rightarrow$  minimize  
- everything else stays the same, this prevents crowding in low-dim space

**Gradient Descent**  
(1) choose  $x_0 \in \mathbb{R}^d$  &  $k=0$   
(2) choose  $\epsilon > 0$  &  $\eta > 0$   
 $x_{k+1} = x_k - \eta \nabla f(x_k)$   
(3) Repeat



**Neural Networks**  
- one hidden layer  
- cross entropy loss:  
 $L = -\sum_{i=1}^n (y_i \ln v_i + (1-y_i) \ln (1-v_i))$   
- ellipsoids with radius  $\sqrt{\lambda}$ , in direction of their corresponding eigenvectors  
- Scaling an MNG is  $\lambda$  matrix by  $X$  scales the isocenters radius by  $\sqrt{\lambda}$   
- CLT:  $\sum_{i=1}^n \text{i.i.d. RVs} \rightarrow$  Gaussian distribution  
- in order for a zero-mean RV  $Z = [Z_1, Z_2]^T$  to be i.i.d.  
 $Z_1, Z_2$  marginally Gaussian  
 $Z_1, Z_2$  are Gaussian  
can be written as  $Z \sim \mathcal{N}(0, \Sigma)$   
two Gaussian RVs can be uncorr. but not indep.  
 $X \sim \mathcal{N}(\mu, \Sigma) \Rightarrow AX + b \sim \mathcal{N}(A\mu + b, A\Sigma A^T)$   
MGF  $(X) = E[\exp(\lambda^T X)] = \exp(\lambda^T \mu + \frac{1}{2} \lambda^T \Sigma \lambda)$   
 $Z = \sqrt{2} (X - \mu) \sim \mathcal{N}(0, I_n)$

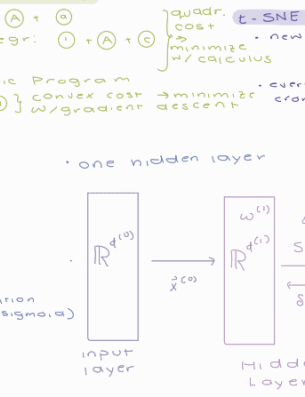
**PCA**

$n$  data points, dimension  $d$   
 $x \in \mathbb{R}^{n \times d}$   
- Procedure:  
(1) Zero-mean the matrix  $X$  to get  $\bar{X}$   
(2) Covariance matrix:  $\Sigma = \bar{X}^T \bar{X}$   
(3) Eigendecomposition, compute  $\Delta = Q \Lambda Q^T$   
(4) keep  $k$  eigenvectors  $Q_k = Q_{:,1:k}$  with most variance (highest eigenvalues in  $\Lambda$ )  
(5) Project your points down to this subspace  
 $\bar{X}_k = \bar{X} Q_k \in \mathbb{R}^{n \times k}$ , these are your principal components  
- if we want our data in the original dimension but only using info retained from PCA, we expand back from PCA to original basis:  
 $\bar{X}$  reconstruct  $= \bar{X}_k Q_k^T \in \mathbb{R}^{n \times d}$   
- Algorithm 2 (SVD)  
(1) Pre-process: reduce mean, normalize variance  
(2) Calculate SVD of  $X: X = U \Sigma V^T$   
(3) New basis is 1st  $k$  right singular vectors (1st  $k$  columns of  $V$ )



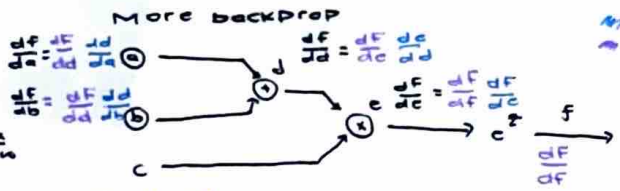
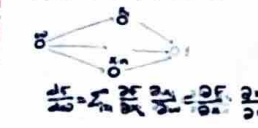
**PCA limitations**  
- assumes only  $\mu$  &  $\Sigma$  affect distribution (ie, Gaussian)  
- due to orthogonality, only exploits linear redundancy  
- isn't robust to scale of features/outliers  
- unsupervised  
**Neighborhood Embedding (NE) / Isomap**  
**Basic NE**  
(1) For each point, find nearest neighbors & draw edge bwn them  
(2) Compute distances bwn points (using only edges)  
(3) Given this set of pairwise distance, put into matrix  $M$  & perform MDS/PCA  
**Stochastic NE**  
(1) Probability  $x_i$  chooses  $x_j$  as its neighbor:  
 $P_{ij} = \frac{\exp\{-\|x_i - x_j\|^2 / 2\sigma^2\}}{\sum_{k:1} \exp\{-\|x_i - x_k\|^2 / 2\sigma^2\}}$   
 $P_{ij} = 0$   
(2) Symmetrize & Normalize  
 $P_{ij} = P_{ji} = \frac{1}{2n} (P_{ij} + P_{ji})$   
 $P = [P_{ij}]$   
(3) Set  $\sigma^2$  adaptively s.t. entropy of  $P_{i\cdot} = \sum_j P_{ij} \log(P_{ij})$   
(4) Posit low-dim representations  $y \in \mathbb{R}^{2k}$  & define stochastic neighborhoods for them  $Q = [Q_i]$   
 $Q_i = \frac{\exp\{-\|y_i - y_j\|^2\}}{\sum_{j:1} \exp\{-\|y_i - y_j\|^2\}}$  (choose  $\sigma = \frac{1}{2}$ )  
Goal: find  $Y$  s.t. BN structure preserved  $\rightarrow$  solve: KL-div  $\text{Fns}$   
 $\hat{y} = \arg \min_Y \text{KL}(P \| Q) = \sum_{i,j} P_{ij} \log \frac{P_{ij}}{Q_{ij}}$   
 $\rightarrow$  use gradient descent to find embedded pts  $\{y_i\}$   
 $\frac{\partial \text{Loss}}{\partial y_i} = \sum_j (P_{ij} - Q_{ij}) C_j = -y_j$   
 $\rightarrow$  nonconvex  $\rightarrow$  many local minima

**Natural Log Rules**  
 $\ln(1) = 0$   $\ln a^b = b \ln a$   
 $\ln(ab) = \ln a + \ln b$   
 $\ln \left(\frac{a}{b}\right) = \ln a - \ln b$   
**Regression Fns**  
(1) Linear:  $v(x; w, a) = wx + a$   
(2) Polynomial: equiv. to linear w/ added poly. feat.  
(3) Logistic:  $v(x; w, a) = \frac{1}{1 + \exp(-wx + a)}$   
**Loss Fns:**  
(A) Squared error:  $L(z, y) = (z - y)^2$   
(B) Absolute error:  $L(z, y) = |z - y|$   
(C) Cross-entropy:  $L(z, y) = -y \ln(z) - (1-y) \ln(1-z)$   
**Cost Fns to minimize:**  
(1) Mean loss:  $J(w) = \frac{1}{n} \sum L(w(x_i), y_i)$   
(2) Max loss:  $J(w) = \max L(w(x_i), y_i)$   
(3) Weighted sum  $J(w) = \sum w_i L(w(x_i), y_i)$   
(4) L2 penalty:  $J(w) = \sum w_i^2 + \lambda \|w\|_2^2$   
(5) L1 penalty:  $J(w) = \sum |w_i| + \lambda \|w\|_1$



**Diagonalizing an Ellipse**  
(1) Compute eigenvalues & eigenvectors of  $A$   
(2)  $Q = [\text{eigenvec}]$ ,  $D = [\text{eigenval}]$   
(3) Change coordinate system along eigenvectors  
 $[x'] = Q^T [x]$   
 $\rightarrow$  Axis points along eigenvectors  $\hookrightarrow$  major & minor axis lengths are  $\sqrt{\lambda_1}$  &  $\sqrt{\lambda_2}$   
 $\Rightarrow$  finding coordinate system axis-aligned  $\equiv$  diagonalizing  $A$

**Backpropagation**



\* = compute during forward pass  
\* = computed during backward pass

**Condition**

Image: 1x3 color channels  $I_1, I_2, I_3$   
 kernel: 3x3 color  $K_1, K_2, K_3$   
 stride: 1x1 color  $S_1, S_2, S_3$   
 convolution operation:  $C_i(x, y) = \sum_{k=1}^3 \sum_{l=1}^3 I_k(x-s_1, y-s_2) \cdot K_{i,l}(x, y)$   
 bias:  $b_i$   
 output:  $O_i(x, y) = C_i(x, y) + b_i$

$$\frac{dL}{dC} = \frac{dL}{dO} \frac{dO}{dC}$$

**Backprop takes  $O(\#edges)$  time**

- convolutions - translationally equivariant (shift input  $\Rightarrow$  shift output fts)

- NN w/ identity fcn as activation fcn is just linear regression

- 2D conv: image  $D \times D$ ,  $N$  filters  $K \times K$ . No padding, output from 1 filter has size  $(D-K+1)(D-K+1)$ . For all  $N$ , mult by  $N$ .

- normalize data in mini-batch  $z^{(i)} = (z^{(i)} - \bar{z}^{(i)}) / \text{Var}(z^{(i)})$
- add scale & shift params  $h^{(i)} = \sigma(\gamma z^{(i)} + \beta)$

# learnable params in NN:

- each image has  $28 \times 28 = 784$  fts
  - 2 hidden layers w/ 100 hidden units each
  - all layers (incl. output) have a bias term
- $\Rightarrow$   $784 \cdot 100 + 100 + 100^2 + 100 + 100$   
 layer 1                      layer 2                      output

For a FC layer  $l$  with  $n_i(l)$  inputs &  $n_o(l)$  outputs,  $W_l$  contains  $n_i(l) \times n_o(l)$  parameters

**Entropy & Information**

entropy: "expected surprise"  
 $H(Y) = E[-\log P(Y)] = -\sum P(Y=x) \log(P(Y=x))$   
 conditional entropy  $H(X|Y) = \sum P(Y=y) H(X|Y=y)$   
 $H(X|Y=y) = -\sum P(X=x|Y=y) \log(P(X=x|Y=y))$   
 information gain / mutual information:  $I(X; Y) = H(Y) - H(Y|X) = H(Y) - \sum P(X=x) H(Y|X=x)$   
 - high entropy: 4y from 4y, uniform dist, sampled values less predictable, flat histogram

**Variational Gradient**

**K-means clustering:**

Partition  $n$  points into  $K$  disjoint clusters

cluster  $i$ 's mean  $\mu_i = \frac{1}{n_i} \sum_{j \in S_i} x_j$

Find  $y$  that min  $\sum_{i=1}^K \sum_{j \in S_i} \|x_j - \mu_i\|^2$

K-means heuristic:

$\rightarrow$  Alt. btwn

- $y$ : fixed; update  $\mu$
- $\mu$ : fixed; update  $y$

$$\hat{W}_{rr} = (X^T X + \lambda I)^{-1} X^T Y$$

$$\hat{W}_{ols} = (X^T X)^{-1} X^T Y$$

$$\hat{\lambda}_{mic} = \frac{1}{n} \sum_{i=1}^n x_i \rightarrow \text{unbiased}$$

$$\hat{\lambda}_{map} = \frac{n}{n+1} \hat{\lambda}_{mic} \rightarrow \text{biased}$$

**Nearest Neighbors**

- non-parametric used for regression & classification  
 - K-NN algo (regression):

Training: store  $\forall (x_i, y_i)$  in dataset  $D$   
 Testing: input:  $x_q$  (query),  $K$  in training dataset

- Find closest  $(\tilde{x}_1, \dots, \tilde{x}_K)$
- set  $\hat{y} = \frac{1}{K} \sum_{j=1}^K f(\tilde{x}_j)$

output:  $\hat{y}$  (prediction)

**Bayes rule:**

$$P(Y|X) = \frac{P(X|Y) P(Y)}{P(X)}$$

$$P(X) = \sum_i P(X|Y=i)$$

$$P(Y|Z, W) = P(Y|Z) P(Y|W)$$

**Classification:**

- Principle: minimize P[error]

$$P[\text{error}] = \int P(\text{error}|x) P(x) dx$$

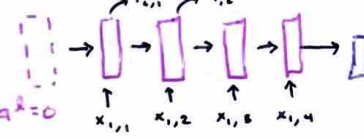
**ResNets:**

- skip connections!

$$x^{(L)} = \sigma(z^{(L-1)}) + x^{(L-1)}$$

(instead of  $x^{(L)} = \sigma(z^{(L-1)})$ )

**RNNs:**



each layer  $a^{l-1}$

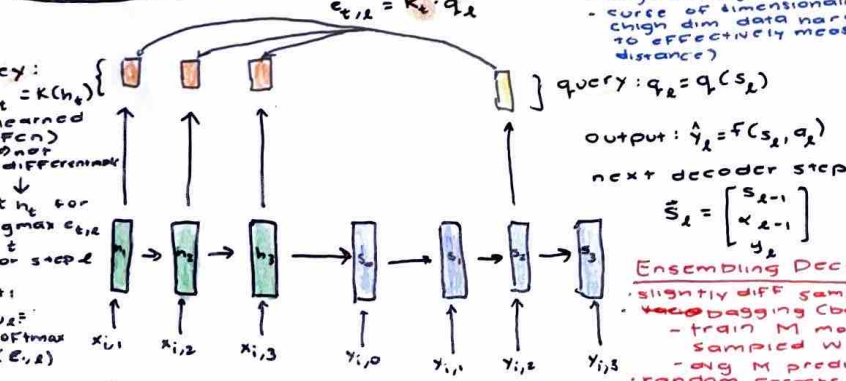
$$a^{l-1} = \begin{bmatrix} x_{i,c} \end{bmatrix}$$

$$z^l = W^l a^{l-1} + b^l$$

$$a^l = \sigma(z^l)$$

$\hat{y} = F(a^L)$   
 auto-regressive gen. model (gen. each new word  $\hat{y}_t$  using fcn in  $\hat{y}_{t-1}$  & previous ones)

**Attention**



Key:  $K_2 = K(h_k)$  (learned fcn) not differentiable  
 set  $h_k$  for argmax  $e_{i,2}$  & for step  $l$   
 Let:  $a_{i,2} = \text{softmax}(e_{i,2})$   
 $e_{i,2} = \frac{\exp(c_{i,2})}{\sum_j \exp(c_{j,2})}$   
 send  $a_{i,2} = \sum_j a_{i,2,j} h_j$

**Ensembling Decision Trees**  
 - slightly diff samples  $\Rightarrow$  diff trees (high variance models)  
 - bagging (bootstrap aggregation)  
 - train  $M$  models, each with  $n'$  (usually  $n/2$ ) samples, sampled w/ replacement  
 - avg  $M$  predictions to get bagged prediction ( $y = \frac{1}{M} \sum G_m(x)$ )  
 - random forests - same as bagging, except at each split, choose only a random subset  $P'$  (usually  $p/2$ ) features to split on  
 - boosting: reweight training points to emphasize ones incor. classified  
 - algo:  
 1. train next model conditioned on all prev models & their weights  
 2. reweight models to minimize loss:  $y = \sum M_i \alpha_i G_i(x)$   
 3. repeat  
 - random forests w/ decision trees of same depth have lower variance than single decision tree

**Decision Trees**

- work well w/ default hyperparams  
 - highly interpretable  
 - NNs generally may perform better  
 - axis-aligned (can't easily model diagonal boundaries)  
 - model may make predictions by posing a series of simple tests on the point  
 - trained in greedy, recursive fashion, downward from root  
 - goal: at each recursion, find the feature  $f$  that split that minimizes conditional entropy (maximizing info. gain)  
 - algorithm for building DTs:  
 1. Start w/ empty tree  
 2. For each node:  
 - IF stopping condition reached:  
 - leaf label = avg. of data at that node  
 - else:  
 - split by next best attribute  
 - argmax  $I(X_i) = \text{argmax } H(Y) - H(Y|X_i)$   
 - recurse to child nodes  
 - stopping conditions:  
 - limited depth: don't split if node if beyond some fixed depth  
 - node purity: don't split if proportion of training point in some class is sufficiently high  
 - info gain criteria: don't split if info gain/purity is sufficiently close to 0  
 - problem: using greedy heuristic (XOR)

- pros of NN:  
 - fast training  
 - learns complex fcn easily  
 - cons of NN:  
 - slow @ test time  
 - high storage cost  
 - curse of dimensionality: high dim data hard to effectively measure distance

- self-attn runtime  $O(d^2 n + n^2 d)$   
 - self-attn & value computation:  
 $K(Q, K, V) = \text{softmax}(QK^T)V$   
 - multi-head attn: have mult. keys, queries & values  
 - self-attn is linear  
 - masked attn typically only used in decoders when computing attn weights @ time  $t$ , weights corr. to times  $t$  w/ set too

**Bias-Variance Tradeoff:** Estimator consists of samples of  $X$ .  
 $X_i$ : new sample  
 $E[(\hat{x} - x)^2] = E[(\hat{x} - \mu)^2] + E[(\mu - x)^2] = \text{Bias}^2 + \text{Var}(\hat{x})$   
 $\text{Bias}^2 = E[(\hat{x} - \mu)^2] = E[(\bar{x} - \mu)^2] = \frac{\sigma^2}{n}$   
 $\text{Var}(\hat{x}) = \text{Var}(\bar{x}) = \frac{\sigma^2}{n}$   
 $E[(\hat{x} - \mu)^2] = \frac{\sigma^2}{n} + \frac{\sigma^2}{n} = \frac{2\sigma^2}{n}$   
**Bias Error** due to inability of hypothesis  $h$  to fit  $D$  perfectly  
**Variance Error** due to fitting random noise in data (usually  $\uparrow$  variance, don't want to add fit unless it reduces bias more)  
**noise in test set** only affects  $\text{Var}(\hat{x})$  but noise in training set also affects bias & variance  
**adding bad fit** rarely increases bias (if set coeffs  $\rightarrow 0$ )

**Markov Models**

$Q = q_1, q_2, \dots, q_n$ : set of  $N$  states  
 $A = a_{11}, a_{12}, \dots, a_{nn}$ : Transition prob matrix  $A$ , each  $a_{ij}$  repr.  $P(q_i \rightarrow q_j)$  s.t.  $\sum_{j=1}^n a_{ij} = 1$   
 $\pi = \pi_1, \dots, \pi_n$ : initial prob distribution over states  
 $\pi_i$ : prob. M.C. will start in state  $i$ .  $\sum_{i=1}^n \pi_i = 1$   
**HMMs** have the above but also have  
 $O = o_1, \dots, o_T$ : sequence of  $T$  observations, each drawn from a vocab  $V = \{v_1, v_2, \dots, v_v\}$   
 $B = b_1, c_1, \dots$ : sequence of observation likelihoods/emission probabilities  $P(o_t | q_t)$  generated from state  $i$ )

**HMM probs:**  
 1) Likelihood: Given specified HMM, compute likelihood of observation sequence  $O$   
 2) Decoding: Given HMM, find best seq. of hidden states

**Probabilistic Graphical Models**

each node repr. random var. & edges repr. dependence relationships  
**DAGs** help us achieve tractability thru cond. independence  
  
**Joint Factorization**  $P(S, X, Y, Z) = P(S, X, Y, Z) = P(S | X, Y, Z) P(X, Y, Z) = P(S | X, Y, Z) P(X | Y, Z) P(Y | Z) P(Z)$   
 each node is conditionally indep. of all of its ancestor nodes, given all of its parents  
 conditional independence  $S \perp\!\!\!\perp X, Y \mid Z \rightarrow P(S, X, Y, Z) = P(S | Z) P(X, Y | Z)$   
 $S \perp\!\!\!\perp Z \mid Y$ ? No.  
 $S \perp\!\!\!\perp X, Y$ ? No.

**Markov Decision Processes**

$S_t \in S$ : state at timestep  $t$ .  $A_t \in A$  action @ time  $t$ .  $\gamma$ : discount factor  
 $R_{t+1} \in R$ : reward generated from  $A_t$ ;  $S_{t+1}$ : state generated from  $A_t$   
**one step dynamics:**  $p(S', r | S, a) = P(S_{t+1}=s', R_{t+1}=r | S_t=s, A_t=a)$   
**return** following time  $t$ :  $G_t = R_{t+1} + \gamma G_{t+1}$   
**policy  $\pi$**  maps states to probability of actions:  $\pi(a|s) = P(A_t=a | S_t=s)$   
**state-value fn**  $V_\pi$  of state  $s$  under policy  $\pi$ :  $V_\pi(s) = E[G_t | S_t=s] = E_\pi[R_{t+1} + \gamma G_{t+1} | S_t=s]$  (expected future return when following  $\pi$  starting at state  $s$ )  
**action-value fn**  $Q_\pi(s, a)$  of taking action  $a$  in state  $s$  under policy  $\pi$ :  $Q_\pi(s, a) = E[G_t | S_t=s, A_t=a] = E_\pi[R_{t+1} + \gamma G_{t+1} | S_t=s, A_t=a]$   
**Bellman expectation eqn:**  
 $V_\pi(s) = \sum_a \pi(a|s) Q_\pi(s, a) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_\pi(s')]$   
**optimal state-value fn** (Bellman optimality eqn)  
 $V_*(s) = \max_a Q_*(s, a) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_*(s')]$   
**Idea:** optimal policy yields max expected total reward  
**Value iteration:** compute optimal values of states by iterative updates until convergence (ie  $V_{k+1}(s) = V_k(s)$ )  
 1)  $V_0(s)$ , initialize  $V_0(s) = 0$   
 2)  $V_k \in S$ , until convergence:  $V_{k+1}(s) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_k(s')]$   
**Policy iteration:** use policy evaluation & extraction to iteratively converge to optimal policy; usually outperforms value iteration bc policies usually converge faster than state values  
 1) Define initial policy  $C$  can be arbitrary  
 2) Until convergence (ie  $\pi_{k+1} = \pi_k$ ):  
 A) Policy evaluation:  $V_{k+1}(s) = \sum_a \pi_k(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_k(s')]$   
 B) Policy improvement:  $\pi'(s) = \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_k(s')]$

**Decision Theory**

**classifier/decision rule  $f$ :**  $X \rightarrow \{1, \dots, K\}$  where  $X$  is space of inputs & we have  $K$  classes  
**loss function  $L(i, j)$ :** penalty classifier receives for predicting class  $i$  when true class is  $j$   
**risk  $R(f)$ :** returns avg loss of classifier for entire data distribution  
 $R(f) = E_{X, Y} [L(f(X), Y)] = E_X [ \sum_{j=1}^K L(f(X), j) P(Y=j | X) ]$   
 $R(f) = \int \sum_{j=1}^K L(f(x), j) P(Y=j | X) P(X) dx = \int \sum_{j=1}^K L(f(x), j) P(X | Y=j) P(Y=j) dx$   
**Bayes classifier:** attains minimum risk a classifier can achieve  
**Bayes risk/error rate:** lowest risk any classifier can attain  
 $f^*(x) = \arg \min_i E_{Y|X} [L(i, Y) | X] = \arg \max_j P(Y=j | X)$   
**equating the risks of 2 classes** should help solve for Bayes classifier decision boundary, eg:  
 $\lambda_{11} P(w_1 | X) + \lambda_{12} P(w_2 | X) = \lambda_{21} P(w_1 | X) + \lambda_{22} P(w_2 | X)$   
 $\lambda_{11} P(w_1 | X) = \lambda_{21} P(w_1 | X) + \lambda_{22} P(w_2 | X) - \lambda_{12} P(w_2 | X)$   
 $\lambda_{11} P(w_1 | X) = \lambda_{21} P(w_1 | X) + (\lambda_{22} - \lambda_{12}) P(w_2 | X)$   
 $\lambda_{11} P(w_1 | X) = \lambda_{21} P(w_1 | X) + \lambda_{21} P(w_2 | X) - \lambda_{12} P(w_2 | X)$   
 $\lambda_{11} P(w_1 | X) = (\lambda_{21} - \lambda_{12}) P(w_2 | X) + \lambda_{21} P(w_1 | X)$   
 $\lambda_{11} P(w_1 | X) - \lambda_{21} P(w_1 | X) = (\lambda_{21} - \lambda_{12}) P(w_2 | X)$   
 $(\lambda_{11} - \lambda_{21}) P(w_1 | X) = (\lambda_{21} - \lambda_{12}) P(w_2 | X)$   
 $\frac{\lambda_{11} - \lambda_{21}}{\lambda_{21} - \lambda_{12}} P(w_1 | X) = P(w_2 | X)$

**MISC**

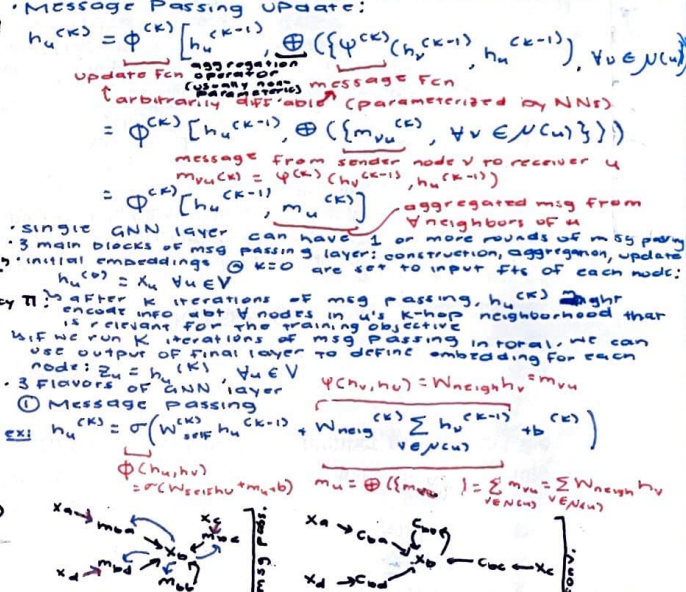
**FC NNs** - more expressive & have greater model capacity due to their incr # of weights (rel. to CNNs)  
**CNNs** less weights & less data need but need  $\uparrow$  architecture params to work well  
**capacity of a model** is a hyperparam that we choose based on the error on a validation set  
 if capacity  $\rightarrow$  optimal  $\rightarrow$  underfitting  
 if capacity  $\rightarrow$  overfitting  
**convolution:** translation equivariant  
**global pooling:** translation invariant  
**self-attn:** permutation equivariant  
**GNN:** permutation invariant bc of the aggregation fn taking as input the representation contained; conv self-attn msg pass bc general msg passing allows arbitrarily irregular graph but conv self-attn needs extra constraints  
**translational equiv.** - useful for pixel-invariant tasks  
**rotational invariance** - useful for graph-invariant tasks  
**tasks:**  
 b CNNs (image-1D): classification/regression on output of conv fn for entire img. eg: dog, car  
 b CNNs (pixel-1D): semantic seg for classification of each individual pixel  
 b GNNs (graph-1D): classification on graphs (e.g. molecular pos/neg. At what temp will it melt?)  
 b GNNs (node/edge 1D): graph of customers & products - deciding pricing of products/recommender systems  
**building decision trees** w/k-ary split; decide k for each node by calculating info gain for diff values of k & optimizing over thresholds & k. Also will prefer high values of  $K$  &  $K$  & a complex decision boundaries & detailed/invariant multi.  
**precision** =  $TP / (TP + FP)$  = TPR  
**recall** =  $TP / (TP + FN)$  = TPR  
**true neg. rate** =  $TN / (TN + FP)$  = 1 - FPR  
 "specificity"  
**deep DT**  $\rightarrow$  model overfit  $\rightarrow$   $\uparrow$  var  
**shallow DT**  $\rightarrow$  model underfit  $\rightarrow$   $\uparrow$  bias

**AUC curves**

**ROC curves:**  
 - Perf. Classifier  
 - False positive rate  
 - FP = (1 - specificity)  
 - TP = sensitivity  
 - AUC (area under curve)  $\rightarrow$  larger area = better model  
 - FN: pred 11, True = 0  
 - FN:  $P = 1, T = 1$   
 - TP:  $P = 1, T = 1$   
 - TN:  $P = 0, T = 0$   
 - ROC curve  $T = 1$  by varying generated by varying threshold of our classifier  
**GNNs**  
 compute  $V$  msg across edge  $\rightarrow$  harder to train, truncate, req. memory  
 might have more param dep. on complexity of parameterization of  $W$   
 grid-like data  $\rightarrow$  conv layer bc of regular neighborhood structure  
 when msg passing a middle interactions bwn neighbors complex only computing scalar-valued averages across edges

**Graph Neural Networks**

nn that can operate on arbitrary relational structures. Able to learn repr. of nodes that depend on structure of the graph  
 message passing used to learn these representations  
 $G = (V, E)$ ; node set  $V$ , edge set  $E$ , set of node feat  $X \in \mathbb{R}^{|V| \times D}$   
 want to generate learned node embeddings  $Z_u$   
 $Z_u, \forall u \in V$   
 during each message passing operation a hidden embedding  $h_u^{(k)}$  is generated, representing the updated embedding of node  $u \in V$  in the  $k$  iteration, based on the info aggregated from its graph neighborhood  $N(u)$  (which can incl.  $u$ )  
**Message Passing Update:**  
 $h_u^{(k)} = \phi \left( \left[ h_u^{(k-1)} \right] \oplus \left( \psi \left( \{ h_v^{(k-1)} \}_{v \in N(u)}, \forall v \in N(u) \right) \right) \right)$   
 update fn aggregation operator (usually non-linear)  
 message fn (parameterized by NNs)  
 message from sender node  $v$  to receiver  $u$   
 $m_{vu}^{(k)} = \psi \left( \{ h_v^{(k-1)}, h_u^{(k-1)} \} \right)$   
 $h_u^{(k)} = \phi \left( \left[ h_u^{(k-1)} \right] \oplus \left( \sum_{v \in N(u)} m_{vu}^{(k)} \right) \right)$   
 aggregated msg from  $v$  neighbors of  $u$   
 3 main blocks of msg passing layer: construction, aggregation, update  
 initial embeddings  $\otimes K=0$  are set to input feats of each node.  
 $h_u^{(0)} = X_u, \forall u \in V$   
 after  $k$  iterations of msg passing,  $h_u^{(k)}$  might encode info abt  $V$  nodes in  $u$ 's  $k$ -hop neighborhood that is relevant for the training objective  
 k: fine run  $k$  iterations of msg passing in total, we can use output of final layer to define embedding for each node:  $Z_u = h_u^{(k)}, \forall u \in V$   
 3 flavors of GNN layer  
 1) Message passing  
 $h_u^{(k)} = \sigma \left( W_{self} h_u^{(k-1)} + W_{neigh} \sum_{v \in N(u)} h_v^{(k-1)} + b \right)$   
 $\psi(h_u, h_v) = W_{self} h_u + W_{neigh} h_v$   
 $m_{vu} = \psi(h_v, h_u) = W_{neigh} h_v + m_{vu}$   
 $m_u = \oplus \{ m_{vu} \} = \sum_{v \in N(u)} W_{neigh} h_v$   
**Note:** if  $W_{self} \in \mathbb{R}^{d \times d}$  the # of learnable params in the MP update fn above is  $2 \cdot d \cdot d$



**Convolutional**  
 $h_u^{(k)} = \phi \left( \left[ h_u^{(k-1)} \right] \oplus \left( \sum_{v \in N(u)} \left\{ W_{self}^{(k)} h_u^{(k-1)} + W_{neigh}^{(k)} h_v^{(k-1)} \right\} \right) \right)$   
 eg: input  $I \in \mathbb{R}^{L \times H \times D}$  (eg image  $L \times H$  w/d  $D$  input chn)  
 $W \in \mathbb{R}^{K \times D \times D}$  (eg filter  $K \times H$  w/d  $D$  in  $L$  &  $H$  input chn)  
 For any point  $(i, j)$  set input embedding bc  
 $h_{ij}^{(k-1)} = [W_{self}^{(k)} h_{ij}^{(k-1)} + \sum_{l, m} W_{neigh}^{(k)} h_{lm}^{(k-1)}]$  (this is value of  $W$  in input chn)  
 $h_{ij}^{(k)} = \sigma \left( \sum_{l, m} W_{neigh}^{(k)} h_{lm}^{(k-1)} + W_{self}^{(k)} h_{ij}^{(k-1)} \right)$   
 \*msg only dep on sender node  $v$  but not  $u$  & weights aren't dep on either  $u$  or  $v$

**Attentional**  
 $h_u^{(k)} = \phi \left( \left[ h_u^{(k-1)} \right] \oplus \left( \alpha \left( h_u^{(k-1)}, h_u^{(k-1)} \right), \left\{ h_v^{(k-1)} \right\}_{v \in N(u)} \right) \right)$   
 $\alpha(h_u, h_v)$  are attn. weights  
 support attn weights are computed by single-head scaled dot-prod  $\alpha = \frac{\exp(\text{dot}(h_u, h_v))}{\sum_{v \in N(u)} \exp(\text{dot}(h_u, h_v))}$   
 is msg passing on FC graph. can parameterize attn weights  
 $\alpha(h_u, h_v) = \exp \left( \frac{1}{\sqrt{d}} W_{attn} h_u \cdot W_{attn} h_v \right)$   
 to emulate self-attn more, could also multiply by  $h_u \cdot h_v$   
 \*weights are  $2 \cdot d \cdot d$  & messages dep on sender & receiver

iterative method to find local ML/MAP estimates of params in statistical models  
coordinate ascent algo

2 steps:

① Expectation  

$$q^{t+1} = \underset{q}{\operatorname{argmax}} F(q, \theta^t) = \underset{q}{\operatorname{argmax}} H(q(z; x; \cdot)) + \mathbb{E}_q [L_c(x; z; \theta^t)]$$

$$= \underset{q}{\operatorname{argmax}} \sum_{k=1}^K q(z_i = k | x_i) \log [q(z_i = k | x_i)] + \sum_{k=1}^K q(z_i = k | x_i) \log p(x_i, z_i = k; \theta)$$

$$q^{t+1}(z_i = k | x_i) = p(z_i = k | x_i; \theta^t) \rightarrow \text{used to compute } \mathbb{E}_q [L_c(x; z; \theta^t)]$$

② Maximization (param estimation)  

$$\theta^{t+1} = \underset{\theta}{\operatorname{argmax}} F(q^{t+1}, \theta) = \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{q^{t+1}} [L_c(x; z; \theta)]$$

EM for MoG

MoG model:  $\tilde{x} | z \sim \mathcal{N}(\mu_z, \Sigma_z)$ ,  $p(z=k) = \alpha_k$   
 $\theta := \{\mu_k, \Sigma_k, \alpha_k\}$ ,  $x_1, \dots, x_n \in \mathbb{R}^d$  observed data  
 $q_{ki}^t := q^t(z_i = k | x_i)$

① E step:  $q^{t+1}(z_i = k | x_i) = \frac{\alpha_k^t p(x_i | z_i = k; \theta^t)}{\sum_{j=1}^K \alpha_j^t p(x_i | z_i = j; \theta^t)}$

② M step:  

$$\mu_k^{t+1} = \frac{\sum_{i=1}^n q_{ki}^{t+1} x_i}{\sum_{i=1}^n q_{ki}^{t+1}} ; \Sigma_k^{t+1} = \frac{\sum_{i=1}^n q_{ki}^{t+1} (x_i - \mu_k^{t+1})(x_i - \mu_k^{t+1})^T}{\sum_{i=1}^n q_{ki}^{t+1}}$$

$$\alpha_k^{t+1} = \frac{1}{N} \sum_{i=1}^n q_{ki}^{t+1}$$

Jensen's inequality:  $f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]$

Distance Fns

- satisfy the following:
  - $d(x, y) = 0$  iff  $x = y$
  - $d(x, y) = d(y, x) \forall x, y$
  - $d(x, z) \leq d(x, y) + d(y, z)$  (triangle ineq.)
- dissimilarity measure  $d(x, y)$  satisfies the above except possibly ③
- similarity measure  $s(x, y) = -d(x, y)$
- computational complexity  $z \in \mathbb{R}^d, D = \{x_1, \dots, x_n\} \in \mathbb{R}^d$
- naive exhaustive search (computes dist b/wn  $z$  &  $\forall$  pts in  $D$  & returns  $k$  nearest points using quickselect)
  - $\hookrightarrow O(d)$ : time computing distances b/wn  $z$  &  $x_i$
  - $\hookrightarrow O(dn)$ : time " " " &  $\forall x_i$
  - $\hookrightarrow O(n)$ : time spent finding  $k$  shortest distances
  - overall:  $O(dn + n) = O(dn)$

divide & conquer: place Cartesian grid over  $\mathbb{R}^d$  space. in this case, assume  $k=1$   
 $\hookrightarrow$  in  $d$  dimensions need to search  $3^d$  cells  
 $\hookrightarrow$  overall  $O(3^d + d3^d) = O(d3^d)$

PCA & L-S

$X \in \mathbb{R}^{n \times d}, Y \in \mathbb{R}^n, \chi = U \Sigma V^T = \sum_{i=1}^d \sigma_i u_i v_i^T$   
 $U \in \mathbb{R}^{n \times d}, \Sigma \in \mathbb{R}^{d \times d}, V \in \mathbb{R}^{d \times d}, \tau \geq \sigma_1 \geq \sigma_2 \geq \dots \geq 0$   
 $\omega_{\text{ridge}} = \sum_{i=1}^d \frac{\sigma_i^2}{\sigma_i^2 + \lambda} v_i u_i^T y$  if  $\lambda \rightarrow 0, \omega_{\text{ridge}} \rightarrow \omega_{\text{ols}}$ , if  $\lambda \rightarrow \infty$  the value of  $\sigma_i$  the less penalized it will be in  $\text{rr}$   
 $\omega_{\text{ols}} = \sum_{i=1}^d \frac{1}{\sigma_i} v_i u_i^T y$   
 $\hookrightarrow$  PCA-OLS =  $\sum_{i=1}^K \frac{1}{\sigma_i} v_i u_i^T y$  multiply by  $\sqrt{K}$  to bring w/PCA to  $d$  dimensions  
 $\hookrightarrow$  large  $\sigma_i$ 's kept intact, while small ones (after  $K$ ) are completely removed - equivalent to  $\lambda \rightarrow \infty$  for 1st  $K$  components &  $\lambda \rightarrow 0$  for rest  
 $\Rightarrow$   $\text{rr}$  is "smooth" version of PCA regr.

Misc

compared to ols,  $\text{rr}$  has  $\uparrow$  bias &  $\downarrow$  variance

## Backpropagation

$f(x_1, \dots, x_n)$ ;  $g: \omega = x$



$$\frac{df}{d\omega} = \sum_{i=1}^n \frac{\partial f}{\partial x_i} \frac{\partial x_i}{\partial \omega} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial \omega}$$

## Convolution

• Image  $I \rightarrow 3$  color channels  $I_r, I_b, I_g$

↳ each channel size  $W \times H$

• Filter  $G \rightarrow 3$  color channels  $G_r, G_b, G_g$

↳ each mask size  $w \times h$

• convolution operation at point  $(x, y)$ :

$$(I * G)[x, y] = \sum_{a=0}^{w-1} \sum_{b=0}^{h-1} \sum_c \{I_c[x+a, y+b] \cdot G_c[a, b]\}$$

↳ size of output of convolution:

• no stride:  $(1+W-w) \times (1+H-h)$

• stride:  $\lfloor 1+(W-w)/s \rfloor \times \lfloor 1+(H-h)/s \rfloor$

$$\frac{\partial L}{\partial G_c[x, y]} = \sum_{i, j} \frac{\partial L}{\partial R[i, j]} \cdot \frac{\partial R[i, j]}{\partial G_c[x, y]}$$

$$= \sum_{i, j} \frac{\partial L}{\partial R[i, j]} \cdot \frac{\partial}{\partial G_c[x, y]} \sum_c \sum_a \sum_b I_c[i+a, j+b] G_c[a, b]$$

$$= \sum_{i, j} \frac{\partial L}{\partial R[i, j]} \cdot I_c[i+x, j+y]$$

• Max pooling; output =  $R$  at  $(i, j)$ :

$$R[i, j] = \max_{a, b} (I * G)[i+a, j+b]$$

## Vanishing Gradient

• when output  $s$  is close to 0 or 1  
 $s' \approx 0 \rightarrow$  GD change  $s$  slowly & the unit is stuck

• mitigation: