

CS 189 Midterm Notes

Multivariate Gaussians: iff each $X_i \in N(\mu_i, \sigma_i^2)$ and each $X_j | X_i \sim N(\mu_{j|i}, \sigma_{j|i}^2)$.
 [or jointly gaussian] iff $P_X(\vec{x}) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-\frac{1}{2}(\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})}$

MLE for Lin Reg: $\vec{w}^* = (X^T X)^{-1} X^T \vec{y}$ is MLE $[\vec{w} | X, \vec{y}]$, where $\vec{y} \sim N(X\vec{w}, \sigma^2)$

MAP for L2 Reg: $\vec{w}^* = (X^T X + \lambda I)^{-1} X^T \vec{y} = \arg \min_{\vec{w}} \|\vec{y} - X\vec{w}\|_2^2 + \lambda \|\vec{w}\|_2^2 = \arg \max_{\vec{w}} \sum_{i=1}^n \log P_{N(\vec{x}_i, \vec{w}, \sigma^2)}(y_i) + \log P_{N(0, \tau^2 I)}(\vec{w})$
 [w/ iid Laplace dist gives L1 Reg] $\vec{w}^* = \arg \min_{\vec{w}} \|\vec{y} - X\vec{w}\|_1 + \lambda \|\vec{w}\|_1 = \arg \max_{\vec{w}} \sum_{i=1}^n \log P_{N(\vec{x}_i, \vec{w}, \sigma^2)}(y_i) + \log P_{N(0, \tau^2 I)}(\vec{w})$
 [Assuming $w_i \sim N(0, \tau^2)$] $\vec{w}^* = \arg \max_{\vec{w}} P(\vec{y} | \vec{w}, X) \cdot P(\vec{w}) = \arg \max_{\vec{w}} P(\vec{w} | X, \vec{y}) = \text{MAP}[\vec{w}]$

Principal Component Analysis: $\Sigma = \text{cov}(X) = X^T X = \text{eig} \Sigma$. Then, columns of Θ_k are k principal components, $X \Theta_k$ = projected data. (k -dim). $D = \text{diag}(\lambda)$. Then, λ_i = % variance explained by i th PC.
 $X \Theta_k$ = d -dim data but only with k -dim information.
 $\text{Tr}(X^T X) = \sum \lambda_i \propto \text{var}(X)$. [assuming zero mean]

Stochastic Neighbour Embeddings: $P_{j|i} = P[X_j \text{ is chosen as } \vec{x}_i \text{'s neighbour}] = \frac{e^{-\|\vec{x}_i - \vec{x}_j\|_2^2 / 2\sigma^2}}{\sum_k e^{-\|\vec{x}_i - \vec{x}_k\|_2^2 / 2\sigma^2}}$, $P_{i|i} = 0$
 For t -SNE, use $\sigma_{j|i} = \frac{\|\vec{y}_i - \vec{y}_j\|_2^2}{\sum_k (\|\vec{y}_i - \vec{y}_k\|_2^2)^2}$ (in place of $\frac{1}{\sigma^2}$, use $\frac{1}{\sigma_{j|i}^2}$).
 This makes tails heavier.

k-means clustering: For each point \vec{x}_i , let $y_i \in [k]$ be the predicted cluster. Then, k -means minimizes $\sum_{i=1}^n \|\vec{x}_i - \vec{\mu}_{y_i}\|_2^2$.
 while !converged: 1) calculate cluster memberships 2) Reset means based on y_i .
 $[y_i = \arg \min_k \|\vec{x}_i - \vec{\mu}_k\|_2^2]$ $[\vec{\mu}_k = \frac{1}{|\{i: y_i = k\}|} \sum_{i: y_i = k} \vec{x}_i]$

Risk: $R(\theta) = E[L(\theta(X), Y)]$.

Logistic Regression: $y_i \in \{0, 1\}$, $\hat{y}_i = \sigma(\vec{w}^T \vec{x}_i)$, $\sigma(z) = \frac{1}{1 + e^{-z}}$. Loss = BCE. Optim-GD or Newton's Method
Binary Cross Entropy: $y \ln \hat{y} + (1-y) \ln(1-\hat{y})$.
 $\vec{w}_{\text{new}} = \vec{w}_{\text{old}} - \eta_k \frac{\partial \mathcal{L}(\vec{w})}{\partial \vec{w}} = \vec{w}_{\text{old}} - \eta_k \frac{\partial \mathcal{L}(\vec{w})}{\partial \vec{w}} \left[\frac{\partial \mathcal{L}(\vec{w})}{\partial \vec{w}} \right]^T$

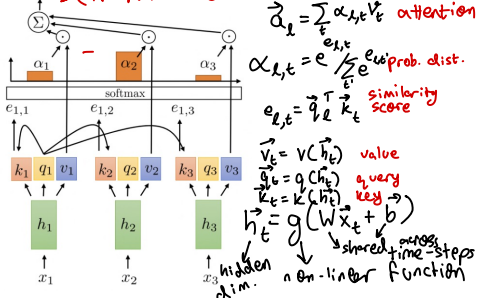
Softmax Regression: $y_i \in [K]$. $\hat{y}_i = \arg \max_j s(\vec{w}_j^T \vec{x}_i)$. $s(z) = \frac{e^{-z_j}}{\sum_k e^{-z_k}}$. Loss = Cross Entropy (negative cross entropy).

Neural Networks: Repeated application of linear map $W^{(l)}$, bias $b^{(l)}$, non-linear activation g .

Back propagation:

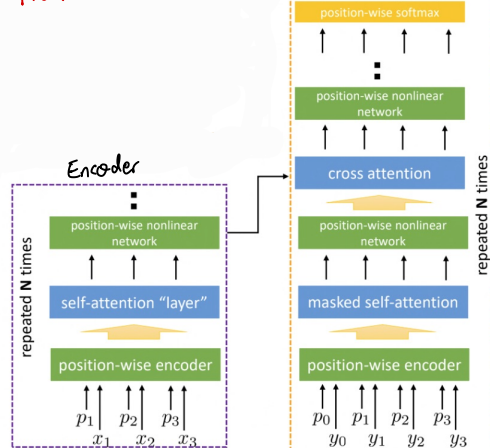
Residual connections:
 - $\vec{z}^{(l)} = \sigma(\vec{z}^{(l-1)}) + \vec{z}^{(l-1)}$
 - smooths out loss landscape
 - allows deeper models
 - ResNet: CNN with RC.

Self-Attention:



Positional Encodings:
 $\vec{p}_k = \begin{bmatrix} \sin(\frac{2\pi k}{10000^{2d}}) \\ \cos(\frac{2\pi k}{10000^{2d}}) \\ \sin(\frac{2\pi k}{10000^{2d-2}}) \\ \cos(\frac{2\pi k}{10000^{2d-2}}) \end{bmatrix}$
 - bounded magnitude
 - non-constant
 - usually uses \sin/\cos .

Transformers:



Batch Normalization:

- Learnable scale and shift
- compute batch statistics (scale, shift) at train time, use to rescale at eval time
- standardization against vanishing gradient.

Masked Self-attention: Send $k_i \rightarrow 0$ if not yet predicted

Derivatives

Cheats:

$$\begin{aligned} -\sigma'(x) &= \sigma(x)(1-\sigma(x)) & -\frac{\partial y^{-1}}{\partial x} &= -y^{-1} \frac{\partial y}{\partial x} y^{-1} \\ -\frac{\partial \|\vec{x}\|_2^2}{\partial \vec{x}} &= 2\vec{x} \\ -\frac{\partial \vec{a}^T X \vec{b}}{\partial X} &= \vec{a} \vec{b}^T & -\frac{\partial \vec{b}^T X^T X \vec{c}}{\partial X} &= X(\vec{b} \vec{c}^T + \vec{c} \vec{b}^T) \end{aligned}$$

$$|\lambda_{\max}(A)| = \|A\|_2$$

Classification

	N	P
N	TN	FP
P	FN	TP

- specificity = $\frac{TN}{TN+FP}$
 - sensitivity = $\frac{TP}{TP+FN}$
 (precision) (recall)

Nearest neighbors

Algorithm The k -nearest neighbors classification algorithm

Input:

D : a set of training samples $\{(x_1, y_1), \dots, (x_n, y_n)\}$
 k : the number of nearest neighbors
 $d(x, y)$: a distance metric
 x : a test sample

1: for each training sample $(x_i, y_i) \in D$ do
 2: Compute $d(x, x_i)$, the distance between x and x_i
 3: Let $N \subseteq D$ be the set of training samples with the k smallest distances $d(x, x_i)$
 4: return the majority label of the samples in N

Decision Trees

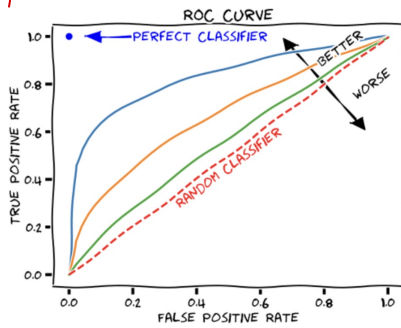
- Entropy, $H(Y) = \mathbb{E}_Y[-\log_2 P(Y)] = -\sum_k P(Y=k) \log_2 P(Y=k) \Rightarrow H(Y|X) = \sum_{x \in X} p(x) H(Y|X=x)$
- Surprise $(Y=k) = -\log_2 P(Y=k)$
- Information Gain: $I(x_{j,v}; Y) := H(Y) - H(Y|X_{j,v})$
- Decision tree nodes split on one attribute based on threshold to maximize Information gain.
 \Rightarrow Axis-aligned decision boundaries. \Rightarrow Recursive splitting, greedy for Info.
- Easily overfit \Rightarrow Ensemble to reduce variance
- Bootstrap Aggregation - ing
 - Train M models with n' (usually $= n$) samples each (w/ replacement) each time
 - Average M predictions $[y = \frac{1}{M} \sum_{m=1}^M G_m(x)]$
- Random Forests: same as Bagging but choose $p' \ll p$ features to split on
- Boosting:
 - Adaptive Boosting:
 - initialize $w_i = \frac{1}{n} \forall i \in [n]$
 - for m in $[M]$:
 - Build model $G_m: \mathbb{R}^d \rightarrow \{-1, 1\}$, where training point x_i is weighted w_i .
 - $e_m = \frac{\sum_i \text{misclassified } w_i}{\sum_i w_i}$. Reweight points: $w_i^* = \begin{cases} \sqrt{\frac{1-e_m}{e_m}} & \text{if misclassified} \\ \sqrt{\frac{e_m}{1-e_m}} & \text{otherwise.} \end{cases}$
 - during inference, $y = \frac{1}{M} \sum_{m=1}^M \alpha_m G_m(x)$. $\alpha_m = \frac{1}{2} \ln(\frac{1-e_m}{e_m})$
 $[e_m \rightarrow 0, \frac{1-e_m}{e_m} = \infty, \alpha_m = \infty; e_m \rightarrow 1, \frac{1-e_m}{e_m} = 0, \alpha_m \rightarrow -\infty]$
 $[$ higher weight to good classifiers $]$.
 - Gradient Boosting: seek m^{th} model output $\vec{g} = \arg \max_{\vec{g} \in G} | \langle -\nabla_{\vec{g}} L(\vec{y}, F_{m-1}(x)), \vec{g} \rangle |$
 where F_{m-1} is averaged output of $m-1$ models so far.

Bias-Variance Tradeoff

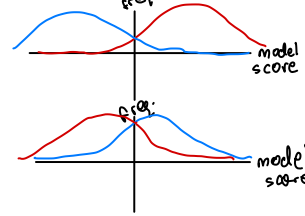
- with true distribution $\vec{y} = \underbrace{f(\vec{x})}_{\text{true relationship}} + \underbrace{z}_{\text{noise}}$, data D , and model h :

$$\mathcal{E}(\vec{x}; h) = \underbrace{(E[h(\vec{x}; D)] - f(\vec{x}))^2}_{\text{model error}} + \underbrace{\text{Var}(h(\vec{x}; D))}_{\text{bias}^2 \text{ of } h} + \underbrace{\text{Var}(z)}_{\text{variance of } h} = |E[\vec{y}] - f(\vec{x})| + \text{Var}(\vec{y}) + \text{Var}(y - f(\vec{x})).$$

irreducible error



- Area Under Curve $\uparrow \Rightarrow$ better model



- cons:

- Slow
- Curse of dimensionality
 [Neighbors far away and hard to find in higher dim. data]

- Pros:

- no training
- non-linear functions can be learned
- 1-NN error $\leq 2 \times$ Bayes error $[\lim_{n \rightarrow \infty}]$
- k-NN error = Bayes error $[\lim_{k, n \rightarrow \infty}]$

Viterbi for HMMs

```
function VITERBI(O, S, Π, Y, A, B): X
  for each state i = 1, 2, ..., K do
    T1[i, 1] ← πi · Biy1
    T2[i, 1] ← 0
  end for
  for each observation j = 2, 3, ..., T do
    for each state i = 1, 2, ..., K do
      T1[i, j] ← maxk (T1[k, j-1] · Aki · Biyj)
      T2[i, j] ← arg maxk (T1[k, j-1] · Aki · Biyj)
    end for
  end for
  zT ← arg maxk (T1[k, T])
  xT ← szT
  for j = T, T-1, ..., 2 do
    zj-1 ← T2[zj, j]
    xj-1 ← szj-1
  end for
  return X
end function
```

Probabilistic Graph Models

- Generalization of Bayes' Nets / Markov Models
- Directed Acyclic Graphs allow tractability.

Markov Decision Processes

- Return: $G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} = R_{t+1} + \gamma G_{t+1}$
- Value: $V_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s]$
- Action-value: $q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a]$
- Policy: $\pi(a|s)$ is Prob. with which we act a from s.
- Bellman Eqs:
 - $\rightarrow V_{\pi}(s) = \mathbb{E}_{\pi}[q_{\pi}(s, a)] \rightarrow q_{\pi}(s, a) = \sum_{s', r} (r + \gamma V_{\pi}(s')) p(s', r | s, a)$
 - $\Rightarrow V_{\pi}(s) = \sum_{a \in A} \pi(a|s) \sum_{s', r} (r + \gamma V_{\pi}(s')) p(s', r | s, a)$

Value Iteration: $V_{k+1}(s) = \max_{a \in A} \sum_{s', r} p(s', r | s, a) [r + \gamma V_k(s')]$

Policy Iteration

- Initialize π and loop to converge:
 - Policy Evaluation: solve Bellman Eqs. for \vec{V}_{π} (value iteration)
 - Policy Improvement: solve for $\pi'(s) = \arg \max_a q_{\pi}(s, a)$

Graph Neural Nets

- General CNNs w/o ordered input
- $\vec{h}_u^{(k)}$ is embedding for node u
- $\psi^{(k)}(\vec{h}_v^{(k-1)}, \vec{h}_u^{(k-1)})$ is a message function that takes in embeddings at layer k-1 for nodes u, v.

- \oplus is some message aggregator. (could add, multiply, vectorize, etc.)
- $N(u)$ are neighbors of u. $\vec{m}_u^{(k)} = \psi^{(k)}(\vec{h}_v^{(k-1)}, \vec{h}_u^{(k-1)})$ is message between u and v.
- $\phi^{(k)}(\vec{h}_u^{(k-1)}, \vec{m}_u^{(k)})$ is some update function enacting layer transitions in GNN.
- General form: $\vec{h}_u^{(k)} = \phi^{(k)}(\vec{h}_u^{(k-1)}, \oplus_{v \in N(u)} \psi^{(k)}(\vec{h}_v^{(k-1)}, \vec{h}_u^{(k-1)}))$

- Example: CNN layer with 3x3 kernel, W, sigmoid activation, σ , and "node" (= coords) (x, y):

$$\vec{h}_u^{(k)} = \sigma\left(\sum_{i,j \in \{-1,0,1\}} W_{ij} \vec{h}_{u+i, y+j}^{(k-1)}\right)$$

- Aggregators must be permutation invariant [permutation of inputs irrelevant].
- Hence, also translational, rotational, perspective, reflective invariance.

Miscellaneous

- ISOMAP is a non-linear dim. reduction method.
- Lasso can be used to remove unwanted features (promotes zeroing out weights).

- t-SNE non-linear + stochastic, PCA linear + deterministic.

- OLS is MLE, Ridge is MAP with gaussian prior on weights.
- Lasso is MAP with laplace prior on weights.

- Decision trees are non-differentiable w.r.t params. solvers

- SVMs are trained using quadratic programming

- SVMs have hinge loss w/ ridge penalty.

- Hard margin SVM is soft-margin SVM with $\lim_{C \rightarrow \infty}$.

- KL minimization in SNE / t-SNE is an MLE objective.

- Decision trees are built greedily for practicality > optimality

- 1 v.s. All classification \leq softmax multiclass for mutually excl.