

Linear Regression

$$h(x) = \theta^T x = \theta_0 + \sum_{i=1}^d \theta_i x_i$$

$$\mathcal{J}(\theta) = \frac{1}{2n} \sum_{i=1}^n [h(x^{(i)}) - y^{(i)}]^2$$

$$\frac{\partial \mathcal{J}}{\partial \theta} = \frac{1}{n} x^T x \theta - x^T y$$

$$\text{GD}$$

$$\theta \leftarrow \theta - \alpha \frac{1}{n} \sum_{i=1}^n [h(x^{(i)}) - y^{(i)}] x^{(i)}$$

$$\bar{\theta} \leftarrow \theta - \alpha \nabla \mathcal{J}(\theta)$$

$$\frac{\partial \mathcal{L}}{\partial \theta} \leftarrow \bar{\theta} - \alpha \sqrt{\nabla \mathcal{L}(x^{(i)}, \theta)} \quad \text{s.t. } l = (h - y)^2$$

$$\bar{\theta} \leftarrow \theta - \alpha \left[h(\bar{x}) - y \right] \bar{x}$$

MGD

$$\theta \leftarrow \theta - \alpha \frac{1}{B} \sum_{i \in B} \nabla \mathcal{L}(x^{(i)}, \theta) \quad \text{s.t. } l = (h - y)^2$$

Norm

$$\hat{\theta} = (x^T x)^{-1} x^T y$$

Linear Basis Gen.

$$h(x) = \theta^T \phi(x) = \sum_j \theta_j \phi_j(x)$$

Loss = MSE (same)

$$\text{GD}$$

$$\theta \leftarrow \theta - \alpha \frac{1}{n} \sum_{i=1}^n [h(x^{(i)}) - y^{(i)}] \phi(x^{(i)})$$

$$\hat{\theta} = [\phi(x)^T \phi(x)]^{-1} \phi(x)^T y$$

Generalization

K-fold CV

- $k = n \Rightarrow$ leave 1 out
- choose best

Ridge Regression

$$\mathcal{J}(\theta) = \frac{1}{2n} \sum_{i=1}^n [h(x^{(i)}) - y^{(i)}]^2 + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2$$

do L_2 regularized

$$\theta \leftarrow \theta - \alpha \frac{1}{n} \sum_{i=1}^n [h(x^{(i)}) - y^{(i)}] x^{(i)}$$

$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n [h(x^{(i)}) - y^{(i)}] x_j^{(i)} - \alpha \lambda \theta_j$$

$$\bar{\theta} = [x^T x + \lambda I_{d \times d}]^{-1} x^T y$$

Perceptrons

$$h(\bar{x}) = \text{sign}(\theta^T \bar{x}) = \begin{cases} 1 & \geq 0 \\ -1 & \leq 0 \end{cases}$$

Surrogate Loss

$$\mathcal{J}(\theta) = \frac{1}{n} \sum_{i=1}^n \max(0, -y^{(i)} \theta^T x^{(i)})$$

$$h(x)^T = \begin{bmatrix} 1 \\ 0 & -1 \\ 2 & 0 \end{bmatrix}$$

$$\bar{\theta} \leftarrow \bar{\theta} - \frac{\alpha}{2} \sum_{i=1}^n [h(x^{(i)}) - y^{(i)}] x^{(i)}$$

Logistic Regression

$$h(\bar{x}) = \sigma(\theta^T \bar{x}) \quad \text{s.t. } \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$

Iteration Loss

$$\ell(\theta) = \begin{cases} -\log[h(\bar{x})] & y^{(i)} = 1 \\ -\log[1 - h(\bar{x})] & y^{(i)} = 0 \end{cases}$$

Candidate Loss

$$\mathcal{J}(\theta) = -\sum_i y^{(i)} \log[h(x^{(i)})] + (1 - y^{(i)}) \log[1 - h(x^{(i)})]$$

$$\text{GD}$$

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_i (h(x^{(i)}) - y^{(i)})$$

$$\theta_j \leftarrow \theta_j - \alpha \sum_i (h(x^{(i)}) - y^{(i)}) x_j^{(i)} + \lambda \theta_j$$

Matrices

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \times \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} = \begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} & a_{11}b_{12} + a_{12}b_{22} \\ a_{21}b_{11} + a_{22}b_{21} & a_{21}b_{12} + a_{22}b_{22} \end{bmatrix}$$

$$A^{-1} A = A^{-1} A = I \Rightarrow \text{if } 2 \times 2 \quad A^{-1} = \frac{1}{|A|} A$$

$$(AB)^T = B^T A^T \quad \& \quad (ABC)^T = C^T B^T A^T$$

Multi-Class Logistic Regression

For 2 classes:

$$h_\theta(x) = \frac{1}{1 + \exp(-\theta^T x)} = \frac{\exp(\theta^T x)}{1 + \exp(\theta^T x)}$$

weight assigned to $y = 0$ weight assigned to $y = 1$

For C classes:

$$h_{\theta_{1:C}}^{(C)}(x) = p_{\theta_{1:C}}(y = c | x) = \frac{\exp(\theta_{1:c}^T x)}{\sum_{k=1}^C \exp(\theta_k^T x)}$$

- Here θ_c is a parameter vector for class $c \in \{1, \dots, C\}$
- Hypothesis also called the softmax function
- Note that sum of class probabilities equals 1

Multi-Class Logistic Regression

The hypothesis for class c

$$h_{\theta_{1:C}}^{(C)}(x) = \frac{\exp(\theta_{1:c}^T x)}{\sum_{k=1}^C \exp(\theta_k^T x)} \Rightarrow C(d+1) = |\theta|$$

Gradient descent simultaneously updates all parameters for c

- Same derivative as before, just with the above $h_{\theta_{1:C}}^{(C)}(x)$

Predict class label as the most probable label

$$\hat{y} = \arg \max_{c \in \{1, \dots, C\}} h_{\theta_{1:C}}^{(C)}(x)$$

K-Nearest Neighbors

$$h(x_i) = \text{majority} \{y_{\text{nn}_1}(x_i), \dots, y_{\text{nn}_K}(x_i)\}$$

Decision Trees

$$\text{nn}(\bar{x}) = \arg \min_{i \in [n]} \sum_{j=1}^d (x_j - \bar{x}_j)^2$$

$$I(A, Y) = H(Y) - H(A|Y) = \arg \min_{i \in [n]} \| \bar{x} - \bar{x}^{(i)} \|_2^2$$

- Symmetric, ≥ 0

$$H(Y) = -\sum_k P(A=a_k) \log P(A=a_k)^2 \text{ prob of pred if pred }$$

$$H(Y|A) = \sum_k P(A=a_k) H(Y|A=a_k) \text{ prob of features given pred if pred }$$

function BuildTree(D, A, Y)

```

    // D: training set, S: input attributes, Y: class attribute
    if D is empty
        return leaf node with failure
    else if all samples in D have same label label for Y
        return leaf node with label
    else if A is empty or all D have same feature values
        return leaf node with majority vote of values of Y in D
    else
        A ← "best" decision attribute among S using D
        tree ← new tree with root node assigned attribute A
        for each value  $a_k$  of A
             $D_k$  ← examples from D with value  $a_k$  for attribute A
            subtree ← BuildTree( $D_k$ ,  $S - A$ , Y)
            add subtree as child of tree with branch labeled  $A = a_k$ 
        return tree
    
```

Representer Theorem

$$\hat{y} = \sum_i \beta_i \phi(\bar{x}^i)$$

$$\beta_i = -\frac{1}{2\lambda n} L'(\theta^T \phi(x), y)$$

Kernel function

$$K(x^i, \bar{x}^j) = \phi(x^i)^T \phi(\bar{x}^j)$$

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n L(\sum_{j=1}^n \beta_j K(x^{(i)}, x^{(j)}), y^{(i)}) + \lambda \sum_{i=1}^n \sum_{j=1}^n \beta_i \beta_j K(x^{(i)}, x^{(j)})$$

- Standard objective
- Gradient Update
- Kernaled Gradient Update
- Kernaled Gradient Update

Renormalized Hypothesis

$$h(\bar{x}) = \sum_i \beta_i K(\bar{x}, \bar{x}) = K(u, u) \pm 2K(u, v)$$

$$\|\phi(\bar{x})\|^2 = \phi(\bar{x})^T \phi(\bar{x}) = K(\bar{x}, \bar{x})$$

$$\|\phi(\bar{x} - \phi(v))\|^2 = \|\phi(\bar{x})\|^2 + \|\phi(v)\|^2 - 2\phi(\bar{x})^T \phi(v)$$

$$\|\phi(\bar{x}) + \phi(v)\|^2 = \|\phi(\bar{x})\|^2 + \|\phi(v)\|^2 + 2\phi(\bar{x})^T \phi(v)$$

Derivatives of Logarithmic Functions

$$\frac{d}{dx} \ln(x) = \frac{1}{x} \quad \frac{d}{dx} |x| = \frac{x - \Delta}{|x| \Delta x}$$

$$\frac{d}{dx} \ln[f(x)] = \frac{1}{f(x)} f'(x)$$

$$\frac{d}{dx} \log_a(x) = \frac{1}{x \ln a}$$

$$\frac{d}{dx} \log_a[f(x)] = \frac{1}{f(x) \ln a} f'(x)$$

RBF Kernel:

$$k(\bar{x}^{(i)}, \bar{x}^{(j)}) = \exp \left[-\frac{\|\bar{x}^{(i)} - \bar{x}^{(j)}\|_2^2}{2\sigma^2} \right]$$

determines L2 significance

Regression

Given training instances, generate a hypothesis function
 $h(x) \approx y$.

1. Hypothesis class \mathcal{H}
2. Loss function $J(\theta)$
3. Optimize loss

1. Hypothesis class \mathcal{H}
2. Loss function $J(\theta)$
3. Optimize loss

Linear Regression

1. Parameterized linear functions:

$$h_\theta(x) = \theta_0 + \sum_{i=1}^d \theta_i x_i$$

2. Least squares: $J(\theta) = \frac{1}{2n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2$
 (Vectorized: $J(\theta) = \frac{1}{2n} (X\theta - y)^t (X\theta - y)$)
3. Solve $\min_\theta J(\theta)$:

- (a) Gradient Descent: $\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{n} \sum_{i=1}^n \left(\sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) x_j^{(i)}$$

- (b) Stochastic Gradient Descent:

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(x^{(i)}, y^{(i)}, \theta)$$

- $\theta = \theta - \alpha \nabla_\theta \ell(x^{(i)}, y^{(i)}, \theta)$

- (c) Mini-batch Gradient Descent
- (d) Normal Equations: Closed form solution:

$$\theta = (X^t X)^{-1} X^t y$$

Generalization

Generalize the hypothesis class to: $h_\theta(x) = \theta^t \phi(x)$, where ϕ is a k -dimensional basis, eg.

$\phi(x) = [1 \ x \ x^2 \ \dots]^t$ (polynomial regression).

New normal equation: $\hat{\theta} = (\phi(X)^t \phi(X))^{-1} \phi(X)^t y$

Regularization

New loss function:

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2} \|\theta\|^2$$

Representor Theorem

If $\lambda > 0$, then the optimal solution $\hat{\theta}$ is a linear combination of the training examples:

$$\hat{\theta} = \sum_{i=1}^n \alpha_i \phi(x^{(i)})$$

Kernel Trick

1. Kernel function: $K(x, z) = \phi(x)^t \phi(z)$
2. Kernelized hypothesis: $h_\theta(x) = \sum_{i=1}^n \alpha_i K(x^{(i)}, x)$
3. Kernelized loss:

$$J(\alpha) = \frac{1}{n} \sum_{i=1}^n L\left(\sum_{j=1}^n \alpha_j K(x^{(i)}, y^{(j)})\right) + \frac{\lambda}{2} \alpha^t K \alpha$$

A kernel is valid if it can be induced via a feature map ϕ .

Mercer's Theorem

A function $K(x, z)$ is a valid kernel iff it is symmetric and positive semi-definite.

Kernel Composition Rules

1. $K(x, z) = aK_1(x, z) + bK_2(x, z)$
2. $K(x, z) = K_1(x, z)K_2(x, z)$
3. $K(x, z) = cK_1(x, z)$
4. $K(x, z) = g(x)K_1(x, z)g(z)$, for g real-valued.
5. $K(x, z) = \exp K_1(x, z)$

Classification

Given labelled training examples, generate a hypothesis function $h(x) \approx y$.

Perceptron

1. $h_\theta(x) = \text{sign}(\theta^t x)$
2. 0/1 loss: $J_{0/1}(\theta) \frac{1}{n} \sum_{i=1}^n \ell_{0/1}(x^{(i)}, y^{(i)}, \theta)$ where

$$\ell_{0/1}(x^{(i)}, y^{(i)}, \theta) = \begin{cases} 0 & h_\theta(x^{(i)}) = y^{(i)} \\ 1 & \text{otherwise} \end{cases}$$

Surrogate loss:

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \max(0, -y^{(i)} \theta^t x^{(i)})$$

3. $\theta = \theta + \alpha y^{(i)} x^{(i)}$ if $y^{(i)} \theta^t x^{(i)} \leq 0$.

Convergence is guaranteed if the data is linearly separable.

Logistic Regression

1. $h_\theta(x) = \frac{1}{1+e^{-\theta^t x}}$
2. Logistic loss: $J(\theta) = \frac{1}{n} \sum_{i=1}^n \ell_{\text{log}}(x^{(i)}, y^{(i)}, \theta)$

where

$$\ell_{\text{log}}(x^{(i)}, y^{(i)}, \theta) = \begin{cases} -\log(h_\theta(x^{(i)})) & y^{(i)} = 1 \\ -\log(1 - h_\theta(x^{(i)})) & y^{(i)} = 0 \end{cases}$$

3. $\theta = \theta - \alpha \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})$

Multi-Class Logistic Regression

1. $h_\theta(x) = \frac{\exp(\theta^t x)}{\sum_{j=1}^k \exp(\theta_j^t x)}$
2. Multiclass logistic loss:

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \ell_{\text{log}}(x^{(i)}, y^{(i)}, \theta) \text{ where}$$

$$\ell_{\text{log}}(x^{(i)}, y^{(i)}, \theta) = -\log \left(\frac{e^{\theta^t x^{(i)}}}{\sum_{j=1}^k e^{\theta_j^t x^{(i)}}} \right)$$

k-Nearest Neighbors

Nearst neighbor: $nn(x) = \operatorname{argmin}_{x^{(i)} \in S} d(x, x^{(i)})$, where d is a distance metric.

k -nearest neighbors is the set of k nearest neighbors of x .

Then we take the majority vote of the labels of the k -nearest neighbors.

$h(x) = \text{majority } knn(x)$

Takes $O(nd)$ time. Memory intensive. No parameters.

Decision Trees

Given a set of labelled training examples, we want to learn a tree. Useful for interpretability.

Use Occam's razor: prefer simpler trees.

Iterative Dichotomiser (ID3)

1. Choose the best attribute to split on.
2. Recurse on the subsets.

How to choose the best attribute? Use maximum information gain.

Entropy: $H(S) = -\sum_{y \in Y} p(y) \log p(y)$

Information gain: $I(A, Y) = H(Y) - H(Y|A)$

Pre-Pruning

Common options: minimum leaf size, maximum depth, maximum number of nodes.

Post-Pruning

Reduced error pruning: remove a node if it doesn't decrease the error on the validation set.

Other

Bias-Variance Tradeoff

Bias: $\mathbb{E}[\hat{f}(x)] - f(x)$

Variance: $\mathbb{E}[(\hat{f}(x) - \mathbb{E}[\hat{f}(x)])^2]$

SVMs

Margin!

$$y^{(i)} = \gamma^{(i)} \frac{\theta^T x^{(i)}}{\|\theta_{1:d}\|_2}$$

$$\gamma = \min_{i \in [n]} \gamma^{(i)}$$

Hard Margin:

$$\min_{\theta} \frac{1}{2} \|\theta\|_2^2 \text{ given } y^{(i)} \theta^T x^{(i)} \geq 1 \quad \forall i$$

Soft Margin:

$$\min_{\theta} \frac{1}{2} \|\theta\|_2^2 + C \sum_{i=1}^n \epsilon_i$$

$$\epsilon_i \geq 0 \quad \forall i \in [n]$$

Hinge Loss:

$$h_\theta(x) = g(\theta^T x) = \text{sign}(\theta^T x)$$

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n L(\theta^T x^{(i)}, y^{(i)}) + \lambda \|\theta_{1:d}\|_2^2$$

$$L = \max(0, 1 - y^{(i)} \theta^T x^{(i)})$$

Recall Regularized Linear models

$$h_\theta(x) = g(\theta^T x)$$

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n L(\theta^T x^{(i)}, y^{(i)}) + \lambda \|\theta_{1:d}\|_2^2$$

where

- Linear Regression: $g(z) = z$ and $L(\theta^T x^{(i)}, y^{(i)}) = (\theta^T x^{(i)} - y^{(i)})^2$
- Perception: $g(z) = \text{sign}(z)$ and $L(\theta^T x^{(i)}, y^{(i)}) = \max(0, -y^{(i)} \theta^T x^{(i)})$
- Logistic Regression: $g(z) = \text{sigmoid}(z)$ and $L(\theta^T x^{(i)}, y^{(i)}) = -y^{(i)} \log(\text{sigmoid}(\theta^T x^{(i)})) - (1 - y^{(i)}) \log(1 - \text{sigmoid}(\theta^T x^{(i)}))$

For SVMs, $g(z) = \text{sign}(z)$. Can we find an $L(\cdot)$ to fit SVMs in this template?

Ensembles

Bootstrap: Create random train sets of n' w/repetition

Bagging: Create many bootstrap replicas and train many classifier on out-of-bootstrapping data average all classifiers for prediction

Random Forests:

AdaBoost:

- Init uniform weight vector \vec{w}_t for n instances s.t. $|\vec{w}_t| = n$
- for $t = 1, \dots, T$:
- Train model $h_t(\cdot)$ on X, y w/ \vec{w}_t
- Measure error ϵ as misclassification
- Choose $\beta_t = \frac{1}{2} \ln \left[\frac{1 - \epsilon_t}{\epsilon_t} \right]$
- Update weight: $w_{t+1,i} = w_{t,i} \exp(-\beta_t y_i h_t(\vec{x}_i)) \quad \forall i$
- Normalize w_{t+1} to be dist. $[0, 1]$ s.t. $\|w_{t+1}\| = 1$

$$w_{t+1,i} = \frac{w_{t,i}}{\sum_j w_{t,j}} \quad \forall i$$

end for

$$\vec{H}(\vec{x}) = \text{sign} \left(\sum_{t=1}^T \beta_t h_t(\vec{x}) \right)$$

if $\epsilon \leq 0.5 \Rightarrow \beta_t \geq 0$ (better than random)

works best on weak learners $h_t \rightarrow$ incorrect one weight inc. logarithmically

K-Means Clustering

Inputs: Training set of n datapoints $(x^{(i)})$, number of clusters k

Step 1: Initialize cluster centroids $\mu_1, \mu_2, \dots, \mu_k$

Step 2: Loop until convergence:

- For every point i , assign it to a cluster $c_i := \arg \min_{j \in [k]} \|x^{(i)} - \mu_j\|_2^2$
- For every cluster j , re-estimate its centroid $\mu_j := \frac{\sum_{i=1}^n \mathbb{1}_{c_i=j} x^{(i)}}{\sum_{i=1}^n \mathbb{1}_{c_i=j}}$

Output: clusters for training points $c_{1:n}$, cluster centers for test points $\mu_{1:k}$

Optimizes $n+d+1$ vars

k-means optimizes the following objective function:

$$\arg \min_S \sum_{i=1}^k \sum_{x \in S_i} \|x - \mu_i\|_2^2$$

where $S = \{S_1, \dots, S_k\}$ is a partitioning over $X = \{x^{(1)}, \dots, x^{(n)}\}$ such that $X = \bigcup_{i=1}^k S_i$ and $\mu_i = \text{mean}(S_i)$

Inertia:

Sum of squared distances to cluster centers

- we want low inertia & low k
↳ but, high $k \rightarrow$ low inertia \rightarrow choose k where $\Delta \text{Inertia} \gg 0$

PCA

Consider an encoding-decoding mechanism for dimensionality reduction:

- Encoding (**Principal Component Analysis or PCA**): $y^{(i)} = Wx^{(i)}$ where $W \in \mathbb{R}^{m \times d}$ and $y^{(i)} \in \mathbb{R}^m$
- Decoding (**Reverse PCA**): $\tilde{x}^{(i)} = Uy^{(i)}$ where $U \in \mathbb{R}^{d \times m}$ and $\tilde{x}^{(i)} \in \mathbb{R}^d$

To learn PCA encoding-decoding, we minimize reconstruction loss:

$$\min_{U, W} \sum_{i=1}^n \|x^{(i)} - \tilde{x}^{(i)}\|_2^2 = \sum_{i=1}^n \|x^{(i)} - UWx^{(i)}\|_2^2 = \|X - XW^T U^T\|_F^2 \text{ (via vectorization)}$$

(where $\|M\|_F^2 = \sum_i \sum_j M_{i,j}^2$ denotes the squared Frobenius norm)

$$\min_{U, W} \sum_{i=1}^n \|x^{(i)} - \tilde{x}^{(i)}\|_2^2 = \|X - XW^T U^T\|_F^2$$

Lemma 1 (without proof): For the PCA objective, the optimal solution for U is such that the columns of U are orthogonal i.e. $U^T U = I_m$ and $W = U^T$

Implication: Can rewrite PCA objective as:

$$\min_U \sum_{i=1}^n \|x^{(i)} - UU^T x^{(i)}\|_2^2$$

s.t. $UU^T = I_d$

Constrained optimization problem

Lemma 2 (without proof): The solution to the PCA learning objective is given by the first m eigenvectors of the empirical covariance matrix

- Let $A = \frac{1}{n-1} \sum_{i=1}^n (x^{(i)} - \mu)(x^{(i)} - \mu)^T$ be the empirical covariance matrix.
- Here, $\mu = \frac{1}{n} \sum_{i=1}^n x^{(i)}$ is the mean vector.
- Let $v_i \in \mathbb{R}^d$ denote the eigenvector for the i -th largest eigenvalue of A . Let $V = [v_1, v_2, \dots, v_d]$ be the matrix of eigenvectors.
- The solution to the PCA learning objective is given by the first m columns of V i.e. $U = [v_1, v_2, \dots, v_m]$

Implication: Can solve for PCA on X by performing an eigendecomposition of A

- v_1, v_2, \dots, v_m are also known as the **principal components directions**
- $v_1, v_2, \dots, v_m x$ are also known as the **principal components scores**

Step 1: Compute covariance matrix $A \in \mathbb{R}^{d \times d}$

$$A = \frac{1}{n-1} (X - \mu)^T (X - \mu) = \frac{1}{n-1} \sum_{i=1}^n (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

Step 2: Compute the eigenvectors of A

- $V, \Lambda \in \text{numpy.linalg.eig}(A)$
- $\{v_i, \lambda_i\}_{i=1 \dots d}$ are the eigenvectors/eigenvalues of A sorted in descending order of eigenvalues

Step 3: Return the top m eigenvectors $U = [v_1, v_2, \dots, v_m]$ as the principal component directions

For any point x , its reduced representation is $y = Wx = U^T x$

To recover: $\vec{y} = U^T \vec{x} \Rightarrow \vec{U} \vec{y} = UU^T \vec{x}$
 $\vec{U} \vec{U}^T = I_d \quad \therefore \quad \vec{x} = U^T \vec{y}$

Explained Variance

measures relative importance of each eigenvector for each principal component i :

$$EV_i = \frac{\lambda_i}{\lambda_1 + \lambda_2 + \dots + \lambda_d} \quad \Rightarrow \text{ignores } d-m \text{ components w/ least significance}$$

Principal Component Direction:

- Order of m -largest eigenvecs show lowest direction of variance
- each is orthogonal to another: $UV^T = I$
- eigen vectors are unit vectors s.t. $\vec{y} = \pm 1/\sqrt{\lambda_i} \vec{v}_i$

Neural Networks

Inputs: $x \in \mathbb{R}^d, y \in \mathbb{R}^c$
 Output: $f_\theta(x) = \text{softmax}(s)$

Logistic Regression

$$s = Wx + b \quad \leftarrow \text{Scores}$$

- Parameters: $W \in \mathbb{R}^{C \times d}, b \in \mathbb{R}^C$

1-hidden layer neural net

$$s = W_2 \max(0, W_1 x + b_1) + b_2$$

- Parameters: $W_1 \in \mathbb{R}^{C \times d}, W_2 \in \mathbb{R}^{C \times H}, b_1 \in \mathbb{R}^H, b_2 \in \mathbb{R}^C$

hidden state is the activated scores

$$s.t. \quad \vec{h} = \max(0, W_1 x + b_1) = \text{ReLU}(\text{scores})$$

Score of next W_2 = $\vec{s} = W_2 \vec{h} + b_2$

Dense Net \equiv Multi-Layer Perception

Activation Functions:

Sigmoid: $\sigma(x) = \frac{1}{1 + e^{-x}}$

tanh: $\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$

ReLU: $\max(0, x)$

Leaky ReLU: $\max(0.2x, x)$

Softplus: $\log(1 + \exp(x))$

ELU: $f(x) = \begin{cases} x & \text{if } x \leq 0 \\ \alpha(\exp(x) - 1) & \text{if } x > 0 \end{cases}$

2-HL NN is a universal approximator like kNN

Forward Pass:

$h = \max(0, W_1 x + b_1)$
 $s = W_2 h + b_2$
 $f_\theta(x) = \text{softmax}(s)$

Cross-Entropy Loss:

$$J(\theta) = - \sum_{i=1}^n \sum_{c=1}^C y_c^{(i)} \log f_\theta(x^{(i)})$$

$$\text{Softmax}(x_i) = \frac{e^{x_i}}{\sum_{c=1}^C e^{x_c}}$$

$$L = - \sum_i y_c^{(i)} \log [\text{softmax}(x_i)]$$

$$\frac{\partial}{\partial x_i} L(\vec{x}^{(i)}, \vec{y}^{(i)}) = s_j^{(i)} - y_j^{(i)} \Rightarrow \frac{1}{N} (\vec{s} - \vec{y})$$

Backprop

$$\frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial f} \frac{\partial f}{\partial \theta} \quad \text{s.t.} \quad \frac{\partial L}{\partial f} = \frac{\partial f}{\partial \theta} \text{ given by back pass}$$

Common Gate reduction:

add gate: gradient distributor

mul gate: "swap multiplier"

copy gate: gradient adder

max gate: gradient router

Vector Derivatives

$x \in \mathbb{R}, y \in \mathbb{R}$
 $x \in \mathbb{R}^N, y \in \mathbb{R}$
 $x \in \mathbb{R}^N, y \in \mathbb{R}^M$

Regular derivative:
 $\frac{\partial y}{\partial x} \in \mathbb{R}$
 $\frac{\partial y}{\partial x} \in \mathbb{R}^{N \times M}$
 $\left(\frac{\partial y}{\partial x} \right)_{i,j} = \frac{\partial y_j}{\partial x_i}$

Derivative is Gradient:
 $\frac{\partial y}{\partial x} \in \mathbb{R}^{N \times M}$
 $\frac{\partial y}{\partial x} \in \mathbb{R}^{M \times N}$

Derivative is Jacobian:
 $\frac{\partial y}{\partial x} \in \mathbb{R}^{M \times N}$

$\nabla_W L(f, y) = \nabla_f L(f, y)^T \vec{H} + \nabla_y \vec{H} = \frac{1}{N} (\vec{s} - \vec{y})^T \vec{H} + \lambda W$
 $\vec{H} = \text{ReLU}(\vec{X}W + b)$ ← activated scores for mat layer

$\nabla_b L(f, y) = \frac{1}{N} (\vec{s} - \vec{y})^T \vec{I}_{0:N}$

Matrix Mults

$y_{i,j} = \sum_k x_{i,k} w_{k,j} \Rightarrow \text{implicit Jacobian}$

$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x} = \left(\frac{\partial x_1}{\partial x_0} \right) \left(\frac{\partial x_2}{\partial x_1} \right) \left(\frac{\partial x_3}{\partial x_2} \right) \left(\frac{\partial L}{\partial x_3} \right)$

$dL/dx = (dL/dy) w^T$
 $[N \times D] \quad [N \times M] \quad [M \times D]$

$dL/dw = x^T (dL/dy)$
 $[D \times M] \quad [D \times N] \quad [N \times M]$

$n_{in} = \text{number of input features}$
 $n_{out} = \text{number of output features}$
 $k = \text{convolution kernel size}$
 $p = \text{convolution padding size}$
 $s = \text{convolution stride size}$

$n_{out} = \frac{n_{in} + 2p - k}{s} + 1$

Dropout: drop nodes w/ some probability in training
 \rightarrow slide filter across image \rightarrow activation maps \rightarrow stack
 \rightarrow backprop on those nodes \rightarrow still used for test

Momentum: moving averages to smooth function S_t

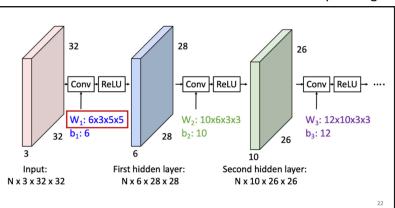
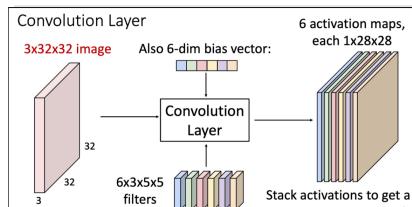
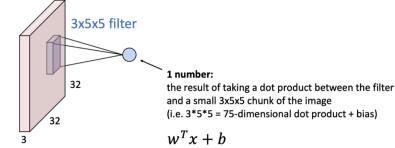
$V_0 = 0$ and $\beta \in [0, 1]$

$V_t = \beta V_{t-1} + (1 - \beta) S_t \Rightarrow S_t = \nabla_\theta J(\theta)$

$\vec{V}_{t+1} \leftarrow \vec{V}_t - \alpha \vec{V}_t : \text{smooths grad. descent}$

Convolutional NN's Visuals

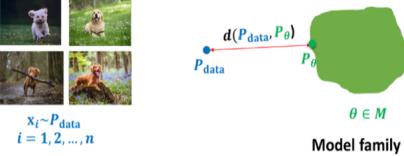
3x32x32 image



Generative Models

Statistical

- Learn a prob. dist. $p_\theta(x) \sim p_{\text{Data}}(x)$
- $p_\theta(x)$ is high when x is deg (e.g.)



Autoregressive Models

- Key idea:** Decompose the joint distribution as a product of 1d probability conditionals using chain rule of probability

$$p_\theta(x) = \prod_{i=1}^d p_\theta(x_i | x_1, x_2, \dots, x_{i-1}) = \prod_{i=1}^d p_\theta(x_i | x_{<i})$$

$$p(x_{1:d}) = \prod_{i=1}^d p(x_i | x_{<i})$$

Output
↑
GPT-2

Representing Prob. Dists.

- Distribution Type:** Depends on whether x_i is discrete or continuous
 - Continuous x_i : p_θ can be a Gaussian distribution
 - Discrete x_i : p_θ can be a categorical (multinomial) distribution
- The parameters θ signify a mapping to the sufficient statistics of the chosen distribution for p_θ
 - Gaussian: θ are parameters of any function that map to mean, variance of p_θ
 - Categorical: θ are parameters of any function that map to probabilities of each category
- Parameters θ can be shared or separate for each $p_\theta(x_i | x_{<i})$
 - Shared: one function with k outputs where k is number of sufficient stats for a distribution type (e.g., $k=2$ for Gaussian, C for categorical)
 - Separated: d functions, each with k outputs

RNNs & Transformers try to parametrize conditionals

Inference requires sequential sampling

$$\begin{aligned} \widehat{x}_1 &\sim p_\theta(x_1) \\ \widehat{x}_2 &\sim p_\theta(x_2 | \widehat{x}_1) \\ \widehat{x}_3 &\sim p_\theta(x_3 | \widehat{x}_1, \widehat{x}_2) \\ &\dots \\ \widehat{x}_d &\sim p_\theta(x_d | \widehat{x}_{<d}) \end{aligned}$$

Eigen Decomposition

$$A\vec{v} = \lambda\vec{v} \rightarrow \text{all } \vec{v} \text{ are eigenvectors}$$

$$(A - \lambda I)\vec{v} = \vec{0} \rightarrow \text{eigenvectors}$$

$$|A - \lambda I| = 0 \rightarrow \text{eigenvectors}$$

Matrix Properties

$$U \text{ is orth.} \Rightarrow UU^T = U^T U = I$$

$$(AB)^T = B^T A^T$$

$$(ABC)^T = C^T B^T A^T$$

$$AB = BA = I \Rightarrow \text{invertible: } B = A^{-1}$$

Max partial

$$\frac{\partial \max(x, x)}{\partial x} = \begin{cases} 1 & x > 0 \\ 0 & x < 0 \end{cases}$$

MAE partial

$$\frac{\partial}{\partial \theta_j} |h_\theta(x) - y| = \frac{h(x) - y}{|h(x) - y|} \frac{\partial}{\partial \theta_j} (h(x) - y) = x_j \cdot \text{sign}(h(x) - y)$$

Example Decision Tree IA

#	Input Attribute				Label
	Electric	Price	Condition	Color	
1	Yes	L	N	R	Yes
2	No	M	N	G	No
3	No	M	U	G	No
4	Yes	H	U	R	Yes
5	Yes	H	N	G	Yes
6	No	L	N	G	No
7	No	H	N	R	Yes
8	Yes	M	U	G	Yes
9	Yes	L	U	G	No
10	No	L	N	G	No

actions: L: Low, M: Medium, H: High, N: New, U: Used, R: Red, G: Green.

$$H(\text{Buy}) = -\frac{5}{10} \log_2 \frac{5}{10} - \frac{5}{10} \log_2 \frac{5}{10} = 1$$

$$H(\text{Buy} | \text{Price} = \text{Low}) = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} = 0.8$$

$$H(\text{Buy} | \text{Price} = \text{Medium}) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 1.6 - \frac{2}{3} = 0.93$$

$$H(\text{Buy} | \text{Price} = \text{High}) = 0$$

$$H(\text{Buy} | \text{Price}) = \frac{4}{10} \times 0.8 + \frac{3}{10} \times 0.93 + \frac{3}{10} \times 0 = 0.6$$

$$I(\text{Buy}, \text{Price}) = 0.4$$

Midterm Kernel Proofs

- (a) [2 points] Let $k(\mathbf{u}, \mathbf{v})$ be a valid kernel function for $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$. Prove or disprove that $ck(\mathbf{u}, \mathbf{v})$ is a valid kernel function for $c > 0$.

Answer: Yes, this is a valid kernel (1pt). The proof can be done in 2 ways:

(1) We can prove it by first principles. Let $\phi(\mathbf{u}, \mathbf{v})$ be the feature map for $k(\mathbf{u}, \mathbf{v})$. Then $\sqrt{c}\phi(\mathbf{u}, \mathbf{v})$ is the feature map for $ck(\mathbf{u}, \mathbf{v})$ (1 pt).

(2) Alternatively, we can prove it by Mercer's Theorem. Symmetry (0.5pts): Since $k(\mathbf{u}, \mathbf{v})$ is a valid kernel map, $ck(\mathbf{u}, \mathbf{v}) = ck(\mathbf{v}, \mathbf{u})$. PSD (0.5pts): Since $k(\mathbf{u}, \mathbf{v})$ is a valid kernel, $\mathbf{z}^T K \mathbf{z} \geq 0$ for all \mathbf{z} . Hence, $\mathbf{z}^T K \mathbf{z} \geq 0$ if $c > 0$.

(b) [3 points] Let $k(\mathbf{u}, \mathbf{v})$ be a valid kernel function for $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ and $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be any real-valued function. Prove that $k(f(\mathbf{u}), f(\mathbf{v}))$ is a valid kernel function.

Answer: Yes, this is a valid kernel. Since $k(\mathbf{u}, \mathbf{v})$ is a valid kernel, it can be expressed as $k(\mathbf{u}, \mathbf{v}) = \phi(\mathbf{u})^T \phi(\mathbf{v})$. Now, consider the feature map $\phi': \mathbf{u} \mapsto f(\mathbf{u})\phi(\mathbf{u})$ and definition of a valid kernel to finish the proof of the claim.

- (c) [5 points] Prove or disprove that $\exp(-\|\mathbf{u} - \mathbf{v}\|^2)$ is a valid kernel function.

[Hint:] (1) You can directly use the result from part b above. (2) Recall from lecture the following composition rule (which can be used without proof): for any valid kernel function $k(\mathbf{u}, \mathbf{v})$, $\exp(k(\mathbf{u}, \mathbf{v}))$ is a valid kernel function.

Answer: Yes, this is a valid kernel. We write $\exp(-\frac{1}{2}\|\mathbf{u} - \mathbf{v}\|^2) = \exp(-\frac{1}{2}(\mathbf{u} - \mathbf{v})^T (\mathbf{u} - \mathbf{v})) = \exp(-\frac{1}{2}\mathbf{u}^T \mathbf{u} + \mathbf{u}^T \mathbf{v} - \frac{1}{2}\mathbf{v}^T \mathbf{v}) = \exp(-\frac{1}{2}\mathbf{u}^T \mathbf{u})\exp(\mathbf{u}^T \mathbf{v})\exp(-\frac{1}{2}\mathbf{v}^T \mathbf{v})$. We then note that $\exp(\mathbf{u}^T \mathbf{v})$ is a valid kernel since $\mathbf{u}^T \mathbf{v}$ is a valid kernel (take $\phi(\mathbf{u}) = \mathbf{u}$) and $\exp(\mathbf{u}^T \mathbf{v})$ is a valid kernel for any valid kernel $k(\mathbf{u}, \mathbf{v})$. We then apply the previous part with $f(\mathbf{u}) = \exp(-\frac{1}{2}\mathbf{u}^T \mathbf{u})$ and $k(\mathbf{u}, \mathbf{v}) = \exp(\mathbf{u}^T \mathbf{v})$. If a student gives a correct alternate solution, give full credit.]

NN additional Grads

$$R = \frac{\lambda}{2} \|w\|_2^2 \quad \text{s.t.}$$

$$\hat{R} = \frac{\lambda}{2} \|(w_1 + \dots + w_n)\|_2^2$$

$$\frac{\partial \hat{R}}{\partial w_1} = \lambda w_1$$

$$S_j = \text{Softmax}(x_j) = \frac{e^{x_j}}{\sum_k e^{x_k}}$$

$$\frac{\partial S_j}{\partial x_j} = S_j \left[1 - \sum_k S_k \right]$$

$$L(\vec{x}, \vec{y}) = -\sum_c^K y_c \log S_c$$

$$L(x, Y) = -\frac{1}{N} \sum_c^K y_c \log S_c$$

$$\frac{\partial L}{\partial x_j} = S_j - y_j$$