

Basics

Fundamental Assumption

Data is iid for unknown P : $(x_i, y_i) \sim P(X, Y)$

True risk and estimated error

True risk: $R(w) = \int P(x, y)(y - w^T x)^2 dx dy = \mathbb{E}_{x,y}[(y - w^T x)^2]$
 Est. error: $\hat{R}_D(w) = \frac{1}{|D|} \sum_{(x,y) \in D} (y - w^T x)^2$
 Generalization Error: Expected - Empirical Standardization

Centered data with unit variance: $\tilde{x}_i = \frac{x_i - \mu}{\sigma}$
 $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$, $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$

Parametric vs. Nonparametric

Parametric: have finite set of parameters. e.g. linear regression, linear perceptron

Nonparametric: grow in complexity with the size of the data, more expressive. e.g. k-NN

Gradient Descent

1. Pick arbitrary $w_0 \in \mathbb{R}^d$
2. $w_{t+1} = w_t - \eta_t \nabla \hat{R}(w_t)$

Stochastic Gradient Descent (SGD)

1. Pick arbitrary $w_0 \in \mathbb{R}^d$
 2. $w_{t+1} = w_t - \eta_t \nabla_{w,t} l(w_t; x', y')$, with u.a.r. data point $(x', y') \in D$ (One / Mini-Batch)
- Momentum: add direction from previous weight update

Regression

Solve $w^* = \operatorname{argmin}_w \hat{R}(w) + \lambda C(w)$

Linear Regression

$$\hat{R}(w) = \sum_{i=1}^n (y_i - w^T x_i)^2 = \|Xw - y\|_2^2$$

$$\nabla_w \hat{R}(w) = -2 \sum_{i=1}^n (y_i - w^T x_i) \cdot x_i$$

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

Ridge regression (=Smaller weights)

$$\hat{R}(w) = \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|_2^2$$

$$\nabla_w \hat{R}(w) = -2 \sum_{i=1}^n (y_i - w^T x_i) \cdot x_i + 2\lambda w$$

$$w^* = (X^T X + \lambda I)^{-1} X^T y$$

L1-regularized regression (Lasso)

$$\hat{R}(w) = \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|_1$$

↳ proximal gradient ↳ not differentiable

Classification

Solve $w^* = \operatorname{argmin}_w l(w; x_i, y_i)$, l ... loss function

0/1 loss (=not convex)

$$l_{0/1}(w; y_i, x_i) = 1 \text{ if } y_i \neq \operatorname{sign}(w^T x_i) \text{ else } 0$$

Perceptron algorithm

Use $l_P(w; y_i, x_i) = \max(0, -y_i w^T x_i)$ and SGD

$$\nabla_w l_P(w; y_i, x_i) = \begin{cases} 0 & \text{if } y_i w^T x_i \geq 0 \\ -y_i x_i & \text{otherwise} \end{cases}$$

Data lin. separable \Rightarrow obtains a lin. separator

Support Vector Machine (SVM)

Hinge loss: $l_H(w; x_i, y_i) = \max(0, 1 - y_i w^T x_i)$

$$\nabla_w l_H(w; y, x) = \begin{cases} 0 & \text{if } y_i w^T x_i \geq 1 \\ -y_i x_i & \text{otherwise} \end{cases}$$

$$w^* = \operatorname{argmin}_w \lambda \|w\|_2^2 + l_H(w; x_i, y_i)$$

For L1-SVM (feature selection) use $\|w\|_1$

Maximize boundary between classes
 Kernels

Properties of kernel

$k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, k must be some inner product (symmetric, positive-definite, linear) for some space \mathcal{V} . i.e. $k(\mathbf{x}, \mathbf{x}') = \langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle_{\mathcal{V}}$
 $\Rightarrow k$ is symmetric and p.s.d.

Kernel matrix

$$K = \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{bmatrix} = \begin{array}{l} \text{Gram matrix} \\ \varphi(\mathbf{x}_1) \cdots \varphi(\mathbf{x}_n) \end{array}$$

Positive semi-definite matrices \Leftrightarrow kernels k

Important kernels

$$\text{Linear: } k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y} = \sum_{i,j} x_i y_j$$

$$\text{Polynomial: } k(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y} + 1)^d$$

$$\text{Gaussian: } k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|_2^2 / (2h^2))$$

$$\text{Laplacian: } k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|_1 / h)$$

Composition rules

Valid kernels k_1, k_2 , also valid kernels: $k_1(x, y) + k_2(x, y)$; $k_1(x, y) \cdot k_2(x, y)$; $c \cdot k_1(x, y)$, $c > 0$;

Reformulating the perceptron

Ansatz: $w^* \in \operatorname{span}(X) \Rightarrow w = \sum_{j=1}^n \alpha_j y_j x_j$

$$\alpha^* = \min_{\alpha \in \mathbb{R}^n} \max_{i=1}^n \left(0, -\sum_{j=1}^n \alpha_j y_i y_j x_j^T x_i \right)$$

Kernelized perceptron and SVM

Use $\alpha^T k$ instead of $w^T x_i$,

use $\alpha^T D_y K D_y \alpha$ instead of $\|w\|_2^2$

$$k_i = [y_1 k(x_i, x_1), \dots, y_n k(x_i, x_n)], D_y = \operatorname{diag}(y)$$

Prediction: $f(\hat{x}) = \operatorname{sign}(\sum_{i=1}^n \alpha_i y_i k(x_i, \hat{x}))$

1. Initialize $\alpha_1 = \dots = \alpha_n = 0$

2. for $t = 1, 2, \dots$, pick data point u.a.r.

3. Predict $\hat{y} = \operatorname{sign}(\sum_{i=1}^n \alpha_i y_i k(x_i, \hat{x}))$

4. if $\hat{y} \neq y \Rightarrow \alpha_i \leftarrow \alpha_i + \eta_i$

Kernelized linear regression (KLR)

Ansatz: $w^* = \sum_{i=1}^n \alpha_i x_i$
 $\alpha^* = \operatorname{argmin}_\alpha \frac{1}{n} \|\alpha^T K - y\|_2^2 + \lambda \alpha^T K \alpha$

$$= (K + \lambda I)^{-1} y$$

Prediction: $f(\hat{x}) = \sum_{i=1}^n \alpha_i k(x_i, \hat{x})$

↳ up/downsampling influence on gradient

Imbalance influence on loss functions
 ↳ it changes the slope

Cost Sensitive Classification

Replace loss by: $l_{CS}(w; x, y) = c_y l(w; x, y)$

Metrics How many selected items are relevant $\rightarrow C_+/C_-$

$$n = n_+ + n_-, n_+ = TP + FN, n_- = TN + FP$$

$$\text{Accuracy: } \frac{TP + TN}{n}, \text{ Precision: } \frac{TP}{TP + FP} = \text{predicted row}$$

$$\text{Recall/TPR: } \frac{TP}{n_+}, \text{ FPR: } \frac{FP}{n_-}$$

$$\text{F1 score: } \frac{2TP}{2TP + FP + FN}$$

$$\text{ROC Curve: } y = \text{TPR}, x = \text{FPR}$$

$$\text{PR curve: } y = \text{precision}, x = \text{recall}$$

Multi-class

↳ for class imbalances

Hinge loss

(Have a weight vector for each class)

$$l_{MC-H}(w^{(1)}, \dots, w^{(c)}; x, y) = \max(0, 1 + \max_{j \in \{1, \dots, y-1, y+1, \dots, c\}} w^{(j)T} x - w^{(y)T} x)$$

↳ 8 class imbalance

- One vs. all: highest confidence

- One vs. one: most number of true

Neural network

↳ slow but simple

c-class / classifiers

Parameterize feature map: $\phi(x, \theta)$ instead of $\phi(x)$, usually: $\phi(x, \theta) = \varphi(\theta^T x) = \varphi(z)$

$$\Rightarrow w^* = \operatorname{argmin}_{w, \theta} \sum_{i=1}^n l(y_i; \sum_{j=1}^m w_j \phi(x_i, \theta_j))$$

Activation functions

$$\varphi'(n_i) = (1 - n_i) \cdot n_i$$

$$\text{Sigmoid: } \frac{1}{1 + \exp(-z)}, \varphi'(z) = (1 - \varphi(z)) \cdot \varphi(z)$$

$$\text{Tanh: } \varphi(z) = \tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$$

$$\text{ReLU: } \varphi(z) = \max(z, 0)$$

Predict: forward propagation

$$v^{(0)} = x; \text{ for } l = 1, \dots, L-1:$$

$$v^{(l)} = \varphi(z^{(l)}), z^{(l)} = W^{(l)} v^{(l-1)}$$

$$f = W^{(L)} v^{(L-1)}$$

Predict f for regression, $\operatorname{sign}(f)$ for class.

Compute gradient: backpropagation

Output layer: $\delta_j = l'_j(f_j)$, $\frac{\partial}{\partial w_{j,i}} = \delta_j v_i$

Hidden layer $l = L-1, \dots, 1$:

$$\delta_j = \varphi'(z_j) \cdot \sum_{i \in \text{Layer}_{l+1}} w_{i,j} \delta_i, \frac{\partial}{\partial w_{j,i}} = \delta_j v_i$$

Learning with momentum

$$a \leftarrow m \cdot a + \eta_t \nabla_W l(W; y, x); W \leftarrow W - a$$

Clustering

k-mean

$$\hat{R}(\mu) = \sum_{i=1}^n \min_{j \in \{1, \dots, k\}} \|x_i - \mu_j\|_2^2$$

$$\hat{\mu} = \operatorname{argmin}_{\mu} \hat{R}(\mu) \quad \text{...non-convex, NP-hard}$$

Algorithm (Lloyd's heuristic): Choose starting centers, assign points to closest center, update centers to mean of each cluster, repeat

- k : prior knowledge / elbow criterion

Dimension reduction

PCA

$$z_i = w^T \cdot x_i \quad x_i = w z_i$$

$$D = x_1, \dots, x_n \subset \mathbb{R}^d, \Sigma = \frac{1}{n} \sum_{i=1}^n x_i x_i^T, \mu = 0$$

$$(W, z_1, \dots, z_n) = \operatorname{argmin} \sum_{i=1}^n \|W z_i - x_i\|_2^2,$$

$$W = (v_1 | \dots | v_k) \in \mathbb{R}^{d \times k}, \text{ orthogonal}; z_i = W^T x_i$$

v_i are the eigen vectors of Σ

Kernel PCA

(Nonlinear dimensionality reduction)

$$\text{Kernel PC: } \alpha^{(1)}, \dots, \alpha^{(k)} \in \mathbb{R}^n, \alpha^{(i)} = \frac{1}{\sqrt{\lambda_i}} v_i,$$

$$K = \sum_{i=1}^n \lambda_i v_i v_i^T, \lambda_1 \geq \dots \geq \lambda_d \geq 0$$

$$\text{New point: } \hat{z} = f(\hat{x}) = \sum_{j=1}^n \alpha_j^{(i)} k(\hat{x}, x_j)$$

Autoencoders

Find identity function: $x \approx f(x; \theta)$

$$f(x; \theta) = f_{\text{decode}}(f_{\text{encode}}(x; \theta_{\text{encode}}); \theta_{\text{decode}})$$

PCA solves relaxed version of k-means problem

line between the centroids for $k=2$ approximates

Probability modeling

Find $h: X \rightarrow Y$ that min. pred. error: $R(h) = \int P(x, y) l(y; h(x)) \partial y \partial x = \mathbb{E}_{x,y}[l(y; h(x))]$

For least squares regression

$$\text{Best } h: h^*(x) = \mathbb{E}[Y | X = x]$$

$$\text{Pred.: } \hat{y} = \hat{\mathbb{E}}[Y | X = \hat{x}] = \int \hat{P}(y | X = \hat{x}) y \partial y$$

Maximum Likelihood Estimation (MLE)

$$\theta^* = \operatorname{argmax} \hat{P}(y_1, \dots, y_n | x_1, \dots, x_n, \theta)$$

$$\theta_{\text{MLE}} = \operatorname{argmin} \sum_{i=1}^n \log p_{\theta}(y_i | x_i)$$

$$\text{E.g. lin. + Gauss: } y_i = w^T x_i + \epsilon_i, \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$\text{i.e. } y_i \sim \mathcal{N}(w^T x_i, \sigma^2), \text{ With MLE (use }$$

$$\operatorname{argmin} -\log l(\theta)$$

$$\theta_{\text{MLE}} = \operatorname{argmin} \sum_{i=1}^n \left[\log p_{\theta}(y_i | x_i) + \sum_{j \neq i} \log p_{\theta}(y_j | x_j) \right]$$

$$\text{Bias/Variance/Noise}$$

$$\text{Prediction error} = \text{Bias}^2 + \text{Variance} + \text{Noise}$$



