

2. DT KNN

Decision Tree

Impurity Measures:

- Misclassification rate: $i_{g(t)} = 1 - \max_c \pi_c$
- Entropy (Shannon): $i_H(t) = -\sum_c \pi_c \log_2 \pi_c$

- Gini index: $i_G(t) = 1 - \sum_c \pi_c^2$

Greedy Optimization: Use $i_H(t)$ or $i_G(t)$ (not $i_{g(t)}$, which doesn't decrease impurity)

Information Gain:

$\cdot I_L = \frac{w_{left}}{N} I_{R_{left}} + \frac{w_{right}}{N} I_{R_{right}}$

$\cdot \Delta I = i(t) - P_L \cdot i(t_L) - P_R \cdot i(t_R)$

Stopping Conditions: $i(t) = 0$ / $dmax / N_{node} < tn$ / $\Delta i(a, t) < \epsilon$

LOOCV: Equivalent to N-fold cross-validation

KNN

Prediction: $\hat{y} = \arg \max_c \sum_{x_i \in N_k(x)} 1(y_i = c)$

Distance Metrics:

- L1 (Manhattan): $d(x_1, x_2) = \sum_i |x_{1,i} - x_{2,i}|$

- L2 (Euclidean): $d(x_1, x_2) = \sqrt{\sum_i |x_{1,i} - x_{2,i}|^2}$

- L ∞ : $d(x_1, x_2) = \max_i |x_{1,i} - x_{2,i}|$

- Cosine Similarity: $\text{Sim}(x_1, x_2) = \frac{x_1^T x_2}{\|x_1\| \|x_2\|}$

- Mahalanobis Distance: $\sqrt{(x_1 - x_2)^T \Sigma^{-1} (x_1 - x_2)}$ (Σ : positive semi-definite, symmetric)

Weighted KNN (inverse distance weighting, closer points more important):

$\cdot \hat{y} = \arg \max_c \frac{1}{n} \sum_{x_i \in N_k(x)} \frac{1}{d(x_i, x)} 1(y_i = c)$

Hyperparameter Selection:

- k small \rightarrow overfitting
- k large \rightarrow underfitting
- Use odd number to avoid ties

Scale Issue: Normalization $x_i = \frac{x_i - \min(x_i)}{\max(x_i)}$ (or use weighted distance)

Confusion Matrix

$^{**}Ground \ Predict^{**}$	$^{**}T^{**}$	$^{**}F^{**}$
$^{**}T^{**}$	TP	FN
$^{**}F^{**}$	FP	TN

Metrics:

- Precision: $\frac{TP}{TP+FP}$
- Sensitivity/Recall: $\frac{TP}{TP+FN}$
- Accuracy: $\frac{TP+TN}{TP+FP+TN}$
- F1 Score: $\frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$

3. Prob Method

Probabilistic Inference

Maximum Likelihood Estimation (MLE):

- $\theta_{MLE} = \arg \max_{\theta} p(D|\theta)$
- $p(D|\theta) = \prod_i p(x_i|\theta)$
- $E_{MLE} = -\ln p(D|\theta) = -\sum_i \ln p(x_i|\theta)$
- $\theta_{MLE} = \frac{\partial \ln p(D|\theta)}{\partial \theta}$

Maximum A Posteriori (MAP):

- $\theta_{MAP} = \arg \max_{\theta} p(\theta|D)$
- $p(\theta|D) \propto p(D|\theta)p(\theta)$
- $E_{MAP} = -(T + a - 1) \ln \theta - (H + b - 1) \ln(1 - \theta)$
- $\theta_{MAP} = \frac{T + a - 1}{T + H + a + b - 2}$

- When $a = b = 1$: $\theta_{MAP} = \theta_{MLE}$

Posterior: $P(\theta|D) = \text{Beta}(\theta|a + T, b + H)$

Hoeffding's Inequality

$p(|\theta_{MLE} - \theta_{true}| \geq \epsilon) \leq 2e^{-2N\epsilon^2} \leq \delta$

Bayesian Models

Predictive Distribution: $p(f|D, a, b) = \int_0^1 p(f|\theta)p(\theta|D, a, b)d\theta$

Fully Bayesian: $\theta^* = \frac{T|T+a|}{T+1+|H+a+b|} = \text{Ber}(f|\theta)$

Conjugate Priors

Bernoulli \Leftarrow Beta $\binom{T}{t}(n-1)!$:

- Likelihood: $p(D|\theta) = \theta^k (1 - \theta)^{n-k}$
- Prior: $\text{Beta}(\theta|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1 - \theta)^{b-1}$

Poisson \Leftarrow Gamma:

- Likelihood: $p(D|\lambda) = \prod_i \frac{\lambda^{x_i} e^{-\lambda}}{x_i!}$
- Prior: $p(\lambda) = \frac{\beta^\alpha}{\Gamma(\alpha)} (\alpha - 1) e^{-\beta \lambda}$
- Posterior: $p(\lambda|D) = \text{Gamma}(\lambda|\alpha + \sum x_i, \beta + N)$

Gaussian \Leftarrow Gaussian:

- Likelihood: $p(D|\mu) = \prod_i \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$
- Prior: $p(\mu) = \mathcal{N}(\mu|\mu_0, \frac{\sigma^2}{2})$
- Posterior parameters:

$\cdot \mu_{\text{post}} = (\mu_0 \sigma_0^2 + N \sigma^2)^{-1} (\mu_0 \sigma_0^2 + \sum x_i \sigma^2)$

$\cdot \sigma_{\text{post}}^2 = 1 / \left(\frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} \right)$

4. Linear Regression

Model: $y_i = f(x_i) + \epsilon_i \quad \epsilon_i \sim \mathcal{N}(0, \beta^{-1})$

Least Squares Error: $E_{LS} = \frac{1}{2} \sum_i (\omega^T x_i - y_i)^2$

Optimal Weight: $\omega^* = \arg \min_{\omega} E_{LS} = (X^T X)^{-1} X^T y$

Non-linear Data (Feature Transform):

- $f(x) = \omega_0 + \sum_{j=1}^M w_j \phi_j(x) = \omega^T \Phi(x)$
- $\Phi \in \mathbb{R}^{N \times (M+1)}$
- $\omega^* = (\Phi^T \Phi)^{-1} \Phi^T y = \Phi_{\dagger}^T y$

Model Complexity:

- High variance \rightarrow overfit
- High bias \rightarrow underfit

Ridge Regression: $E_{ridge} = \frac{1}{2} \sum_i (\omega^T \Phi(x_i) - y_i)^2 + \frac{\lambda}{2} \|\omega\|_2^2$

Probabilistic Formulation

Likelihood: $y_i \sim \mathcal{N}(f_{\omega}(x_i), \beta^{-1})$, $p(y_i|X, \omega, \beta) = \prod_i \mathcal{N}(y_i|f_{\omega}(x_i), \beta)$

Negative Log-Likelihood:

$\cdot E_{ML} = -\ln p(y|X, \omega, \beta)$

$\cdot E_{ML} = \frac{1}{2} \sum_i \mathcal{N}(y_i^T \Phi(x_i) - y_i)^2 + \frac{\lambda}{2} \ln \beta + \frac{N}{2} \ln 2\pi$

Maximum Likelihood Estimators:

$\cdot \omega_{ML} = w_{LS} = \Phi_{\dagger}^T y$

$\cdot \beta_{ML} = \frac{1}{N} \sum_i \Phi(x_i)^T (\omega_{ML} \Phi(x_i) - y_i)^2$

With Gaussian Prior: $p(\omega|\alpha) = \mathcal{N}(\omega|0, \alpha^{-1}I) = \left(\frac{\alpha}{2\pi}\right)^{\frac{M}{2}} \exp\left(-\frac{\alpha}{2} \omega^T \omega\right)$ (M: length of ω)

MAP Estimation:

$\cdot E_{MAP} = -\ln p(y|X, \omega, \beta) - \ln p(\omega|\alpha)$

$\cdot E_{MAP} = \frac{\alpha}{2} \sum_i \mathcal{N}(y_i^T \Phi(x_i) - y_i)^2 + \frac{\alpha}{2} \|\omega\|_2^2$

- Equivalent to Ridge Regression where $\lambda = \frac{\alpha}{\beta}$

$\cdot \omega_{ridge} = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$

Fully Bayesian Linear Regression

Posterior: $p(\omega|D) = \mathcal{N}(\omega|\mu, \Sigma)$

$\cdot \mu = \beta \Sigma \Phi^T y$

$\cdot \Sigma^{-1} = \alpha I + \beta \Phi^T \Phi$

Predictions:

$\cdot \text{MLE: } p(y_{new}|x_{new}, w_{ML}, \beta_{ML}) = \mathcal{N}(y_{new}|\mu_{ML}^T \Phi(x_{new}), \beta_{ML}^{-1})$

$\cdot \text{MAP: } p(y_{new}|x_{new}, w_{MAP}, \beta) = \mathcal{N}(y_{new}|\mu_{MAP}^T \Phi(x_{new}), \beta^{-1})$

• Fully Bayesian:

$p(y_{new}|x_{new}, D) = \mathcal{N}(y_{new}|\mu^T \Phi(x_{new}), \beta^{-1} + \Phi(x_{new})^T \Sigma \Phi(x_{new}))$

Weighted Linear Regression

Objective (with weight γ): $J_{\gamma} = \sum_i \gamma_i (w^T \Phi(x_i) - y_i)^2$

Optimal Weight: $\omega_{\text{weighted}} = (\Phi^T R \Phi)^{-1} \Phi^T R y$

5. Linear Classification

Zero-one Loss: $\ell_0(y, \hat{y}) = \sum_i 1(y_i \neq \hat{y}_i)$ (loss for incorrect predictions is 1)

Hyperplane: $f(x) = w^T x + w_0$

• Direction: w

• Distance from origin: $\frac{|w_0|}{\|w\|}$

Perceptron Update Rule (for each misclassified x_i):

$w \leftarrow \begin{cases} w + x_i & \text{if } y_i = -1 \\ w - x_i & \text{if } y_i = 1 \end{cases} \quad w_0 \leftarrow \begin{cases} w_0 + 1 & \text{if } y_i = -1 \\ w_0 - 1 & \text{if } y_i = 1 \end{cases}$

Probabilistic Generative Model:

• Prior: $y \sim \text{Categorical}(\theta)$, $p(y = c) = \theta_c = \frac{e^{\theta_c}}{\sum_c e^{\theta_c}}$

• Class-conditional: $p(x|y = c) = \mathcal{N}(x|\mu_c, \Sigma)$ (assume Σ_c all equal)

Probabilistic Generative Models & Discriminant Analysis

Binary Classification:

$\cdot p(y = 1|x) = \sigma(a) = \frac{1}{1 + e^{-a}} \text{ where } a = \ln \frac{p(x|y=1)p(y=1)}{p(x|y=0)p(y=0)}$

$\cdot a = w^T x + w_0$

LDA (Linear Discriminant Analysis) (with shared covariance Σ):

$\cdot w = \Sigma^{-1}(\mu_1 - \mu_0)$

$\cdot w_0 = -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 + \ln \frac{p(y=1)}{p(y=0)}$

• Thus $y|x \sim \text{Bernoulli}(\sigma(w^T x + w_0))$

Multi-class Classification:

$\cdot p(y = c|x) = \frac{p(x|y=c)p(y=c)}{\sum_j p(x|y=c_j)p(y=c_j)} = \frac{\exp(w_c^T x + w_{c0})}{\sum_j \exp(w_j^T x + w_{j0})}$

$\cdot w_c = \Sigma^{-1} \mu_c$

$\cdot w_{c0} = -\frac{1}{2} \mu_c^T \Sigma^{-1} \mu_c + \ln p(y = c)$

QDA (Quadratic Discriminant Analysis) (different covariances Σ_c):

$\cdot p(y = 1|x) = \sigma(a) \text{ where } a = x^T W_{2e} + w_1^T x + w_0$

$\cdot W_2 = \frac{1}{2} w_0^{-1} - \Sigma_1^{-1}$

$\cdot w_1 = \Sigma_1^{-1} \mu_1 - \Sigma_0^{-1} \mu_0$

$\cdot w_0 = -\frac{1}{2} \mu_1^T \Sigma_1^{-1} \mu_1 + \frac{1}{2} \mu_0^T \Sigma_0^{-1} \mu_0 + \ln \frac{p_1}{p_0} + \frac{1}{2} \ln \frac{|\Sigma_1|}{|\Sigma_0|}$

Linear Discriminant Model: Logistic Regression

Binary Logistic Regression:

$\cdot p(y = 1|x) = \sigma(w^T x)$

$\cdot p(y = 0|x) = 1 - \sigma(w^T x)$

$\cdot p(y|\omega, x) = \prod_i \mathcal{N}(y_i | \sigma(w^T x_i), (1 - \sigma(w^T x_i))^{1 - y_i})$

Loss Function (Binary Cross Entropy):

$\cdot E(w) = -\sum_i y_i \log \sigma(w^T x_i) + (1 - y_i) \log(1 - \sigma(w^T x_i))$

• Regularization can be added: $+\lambda \|w\|_2^2$

Multi-class (Softmax + Cross Entropy):

$\cdot E(w) = -\sum_i \sum_c y_{ic} \log \frac{\exp(w_c^T x)}{\sum_j \exp(w_j^T x)}$

$\cdot y_{ic} = 1$ iff sample $x \in c$ class

6. Optimization

Convexity

A function is convex if:

1. $f((1-t)x + ty) \leq (1-t)f(x) + tf(y)$ (any point between two points is lower than the line connecting them)
2. $f(t) - f(x) \geq \frac{f(1-t) - f(x)}{1-t-x} (t-x)$
3. $f(t) \geq f(x) + (y-x)^T \nabla f(x)$

4. Hessian Matrix is positive semi-definite

Gradient Descent (Line Search)

1. $\Delta \theta = -\nabla f(\theta)$

2. $t^* = \arg \min_{t \geq 0} f(\theta + t \Delta \theta)$

3. $\theta = \theta + t^* \Delta \theta$

SGD (Stochastic Gradient Descent)

$\theta = \theta + r \cdot \nabla f(\theta)$, where r is learning rate

Decaying learning rate: $r = \alpha r, \quad 0 < \alpha < 1$

Momentum

$\cdot m_t = r \cdot \nabla f(\theta_t) + \gamma \cdot m_{t-1}$

$\cdot \theta_{t+1} = \theta_t - m_t$

Adam

$\cdot m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla f(\theta_t) \text{ (mean)}$

$\cdot v_t = \beta_2 v_{t-1} + (1 - \beta_2) (\nabla f(\theta_t))^2 \text{ (variance)}$

$\cdot \hat{m}_t = \frac{m_t}{1 - \beta_1^t} \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}$

$\cdot \theta_{t+1} = \theta_t - \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}$

• Default values: $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$

Newton Method

Taylor expansion: $f(\theta_t + \delta) = f(\theta_t) + \delta^T \nabla f(\theta_t) + \frac{1}{2} \delta^T \nabla^2 f(\theta_t) \delta + \dots$

Update: $\theta_{t+1} = \theta_t - [\nabla^2 f(\theta_t)]^{-1} \nabla f(\theta_t)$

Mini-batch SGD

$\cdot \theta_{t+1} = \theta_t - r \cdot \frac{1}{|S|} \sum_{j \in S} \nabla L_j(\theta_t)$

• Batch size $\uparrow \rightarrow$ variance \uparrow

• Batch size $\uparrow \rightarrow$ computation time \uparrow

7. Deep Learning

Notation: w_{ijk} denotes weight from layer i , input node j , output node k

Architecture Types

- Feed-Forward Neural Network (FFNN)
- Multi-layered Perceptron (MLP)

Activation Functions

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

• ReLU: $\max(0, x)$

• ELU: $\begin{cases} x & x > 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$

• tanh: $\tanh(x)$

• Leaky ReLU: $\max(0, 1.1x, x)$

• Swish: $x \cdot \sigma(x)$

Target and Loss Functions

Target $Sp(y)$	$x f$	Final Layer	Loss
-----	-----	-----	-----
Binary (Bernoulli)	Sigmoid	BCE (Binary Cross Entropy)	
Discrete (Categorical)	Softmax	CE (Cross Entropy)	

8. CNN & Deep Learning Architecture

CNN Kernel Parameters: $L \times M \times C_{in} \times C_{out}$

Padding

- VALID:** No padding, $D_{l+1} = (D_l - K) + 1$
- SAME:** Add $p = \frac{K-1}{2}$ padding on each side to keep size
- FULL:** Add $p = K - 1$ on each side to increase size

Stride & Pooling

Stride: Step size (< 1 results in downsampling)

$$D_{l+1} = \lfloor \frac{D_l \cdot p - K}{S} + 1 \rfloor$$

Pooling:

- Max pooling:** Take maximum value

- Mean pooling:** Take average value

Initialization & Training Issues

Gradient Problems:

- Vanishing gradient:** W becomes too small

- Exploding gradient:** W becomes too large

Xavier Initialization:

$$W_{fan_in, fan_out} = \frac{2}{\sqrt{fan_in + fan_out}}$$

$$\text{Uniform: } W \sim \text{Uniform}\left(-\sqrt{\frac{6}{fan_in + fan_out}}, \sqrt{\frac{6}{fan_in + fan_out}}\right)$$

$$\text{Normal: } W \sim \mathcal{N}\left(0, \sqrt{\frac{6}{fan_in + fan_out}}\right)$$

- Used for saturating activations like sigmoid and tanh

Regularization & Normalization

- Regularization techniques:**

- Adding ℓ_2 norm (Weight Decay).

- Early stopping.

- Data augmentation.

- Injecting noise.

- Dropout: Used only during training.

Batch Normalization

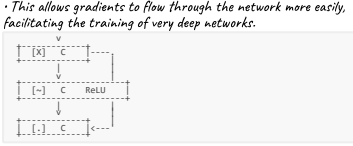
- Standardizes inputs to a layer for each mini-batch:

$$z = \frac{x - \mu_B(x)}{\sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu_B(x))^2}} \Rightarrow z \sim \mathcal{N}(0, 1)$$

Residual Learning (Skip Connections)

- Skip Connection formula:** $y = f(x, W) + x$ (if x is $f(x, W)$)

- This allows gradients to flow through the network more easily, facilitating the training of very deep networks.



9. Support Vector Machines (SVM)

Margin: $\frac{2}{\|w\|}$

Constraints:

- $w^T x_i + b = 1$ for $y_i = +1$

- $w^T x_i + b \leq -1$ for $y_i = -1$

- $\text{Trust: } y_i(w^T x_i + b) - 1 \geq 0$ for $\forall x_i$

Optimization Problem:

- Minimize:** $\frac{1}{2} \|w\|^2$

- Subject to:** $f_i(w, b) = y_i(w^T x_i + b) - 1 \geq 0$

Lagrangian Dual Function

- Dual function:** $g(\alpha) = \min_{w, b \in \mathbb{R}^d} (f_0(w, b) + \sum_{i=1}^M \alpha_i f_i(w, b))$

- Lagrangian:** $L(w, b, \alpha) = f_0(w, b) + \sum_{i=1}^M \alpha_i f_i(w, b)$

- Conditions:** $\alpha_i \geq 0$ and $f_i(w, b) \leq 0$

SVM Optimization Steps

1. Calculate Lagrangian:

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^M \alpha_i [y_i(w^T x_i + b) - 1]$$

2. Minimize L :

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^M \alpha_i y_i x_i \stackrel{!}{=} 0$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^M \alpha_i y_i \stackrel{!}{=} 0$$

$$\Rightarrow w = \sum_{i=1}^M \alpha_i y_i x_i$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

$$\Rightarrow \sum_{i=1}^M \alpha_i y_i x_i = 0$$

Common Examples:

- Polynomial:** $k(a, b) = (a^T b)^n$ or $(a^T b + 1)^P$

- Gaussian (RBF):** $k(a, b) = \exp\left(-\frac{\|a - b\|^2}{2\sigma^2}\right)$

Multiclass Classification Strategies

- 1 vs n classification: Look at maximum distance

- 1 vs 1 classification: Look at majority vote

10. Dimension Reduction PCA SVD

Dimension Reduction (PCA)

$$\text{Transformation: } z' = \bar{x} - \bar{x} \cdot \Sigma_x^{-1} \cdot \bar{x} + F^T \Sigma_x F$$

1. Centering Data

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

$$x_i = x_i - \bar{x} \text{ where } \bar{x} = \frac{1}{N} \sum_{j=1}^N x_j$$

