

**Linear regression** minimize  $\sum \text{mean-square error}$   
**Error:**  $\hat{R}(w) = \|Xw - y\|_2^2$  Assuming  $x^T x$  invertible  
**closed form:**  $w^* = (X^T X)^{-1} X^T y$  full rank  $n > p$   
**Gradient:**  $\nabla_w = 2X^T(Xw - y)$   
**Ridge Regression**  
**Error:**  $\hat{R}(w) = \sum_i (y_i - w^T x_i)^2 + \lambda \|w\|_2^2$   
**closed form:**  $w^* = (X^T X + \lambda I)^{-1} X^T y$   
**Gradient:**  $\nabla_w = -2 \sum_i (y_i - w^T x_i) x_i + 2 \lambda w$   
**Lasso**  $L_1$  penalty penalize  $\beta$  gradient  
 $\beta = \sum_i (y_i - x_i^T \beta)^2 + \lambda \sum_i |\beta_i|$   
 $= (y - X\beta)^T (y - X\beta) + \lambda \|\beta\|_1$   
 No closed form solution.  
**Bayesian Linear Regression**  
 Set a prior over  $\beta$   
 Ridge Regression is equivalent as seeing a Gaussian  
 Prior  $\pi(\beta) \sim N(\beta_0, \frac{\sigma^2}{\lambda}) \propto \exp\{-\frac{\lambda}{2\sigma^2} \beta^T \beta\}$   
 Finding posterior  
**Properties** Involving linear Regression  
 Ridge  $\rightarrow$  smallest var among unbiased  
 biased Unbiased  
 Consistent Consistent  
**Gaussian Process Regression**  
 Rewrite joint dist.  
 $P\begin{bmatrix} y \\ y_* \end{bmatrix} = N\left(\begin{bmatrix} y \\ y_* \end{bmatrix}, \begin{bmatrix} C_n & k \\ k^T & C \end{bmatrix}\right)$   
 $C_n = K + \sigma^2 I_n$   
 $C = K(X, X_{**}) + \sigma^2 I_n$   
 $k = K(x, x_*)$   
 $K = K(X, X)$   
 with  $k$  being the kernel function  
**Prediction with GP**  
 $p(y_* | x_*, X, y) = N(y_* | \mu_*, \sigma_*^2)$   
 $\mu_* = k^T C_n^{-1} y$   
 $\sigma_*^2 = C - k^T C_n^{-1} k$   
 predict new values from  
 $p(y_* | x_*, X, y) = N(y_* | \mu_*, \sigma_*^2)$   
 $\mu_* = k^T (K + \sigma^2 I)^{-1} y$   
**Kernel:** Can map to infinite dimension  
 - computationally cheaper  
 If  $K_1, K_2$  are valid:  $\bigoplus$  valid  $\bigoplus$  valid  $\bigoplus$  valid  $\bigoplus$  valid  $\bigoplus$  valid  $\bigoplus$  valid  
 $k: \mathbb{R}^n \rightarrow \mathbb{R}^d$   $k(K)$  when  $k$  is poly/loop  
 Polynomial kernel  $(x^T x' + c)^d$  touch  $(x^T x' + c)$  is not  
 Radial basis  $\exp(-\frac{\|x - x'\|_2^2}{2\sigma^2})$   $[-1, 1]$   
 To dis- prove P.S.d.  
 find  $f(x)$  s.t.  $\det[K] < 0$   
**Logistics**  $P(w) = \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$  (not)  
 $P(y | x) = \frac{1}{1 + \exp(-y w^T x)}$

**MAP**  
 asymptotic  $P(CND) = \arg \max P(CND | w) \pi(w)$   
 $MLE = MP$  when  $n \rightarrow \infty$  or the prior dist. is uniform  
 Laplace  $-L_2$   
 Gamma  $-L_1$   
 Bias-Variance tradeoff  
 $\text{Bias}(\hat{\beta}) = E[\hat{\beta}] - \beta$   
 $\text{Var}(\hat{\beta}) = E[(\hat{\beta} - E[\hat{\beta}])^2]$   
**Square Error Decomposition**  
 $E_D E_X [(f(x) - Y)^2] = \text{noise}^2$   
 $E_X [(E_Y [Y | X] - Y)^2] = \text{noise}^2$   
 $E_X [E_Y [(f(x) - E_Y [Y | X])^2]] = \text{bias}^2$   
 $E_X [(E_Y [f(x) - E_Y [Y | X]])^2] = \text{bias}^2$   
 Var  $\uparrow$  overfit bias  $\uparrow$  underfit  
 Fit to every pt. All is const.  
**Cross validation / LOO**  $\rightarrow$  unbiased prediction  
 The value of a sample in bootstrap  $(1 - \frac{1}{n})^n \approx 0.67$   
 Therefore, ERM on  $Z$  would have 67% accuracy by memorization  $\rightarrow$  over-confident  
**Bootstrap** type of replacement  
 The value of a sample in bootstrap  $(1 - \frac{1}{n})^n \approx 0.67$   
 Therefore, ERM on  $Z$  would have 67% accuracy by memorization  $\rightarrow$  over-confident  
**Jackknife** A way to de-bias estimator  
 Estimate of an  $\hat{\theta}_n$ 's bias  $\hat{\theta}_n - \theta$   
 $\hat{\theta}_n^{JK} = \hat{\theta}_n - \text{bias}$  JK Estimator (去掉 bias)  
 $\text{bias}^{JK} = (n-1)(\hat{\theta}_n - \hat{\theta}_{n-1})$   
 $\hat{\theta}_n^{JK} = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_{n-1}^{(i)}$  is the avg. LOO Estimator  
 It works for any distribution.  
**Bootstrap debiased:**  $\hat{\theta} = 2\hat{\theta}_n - \frac{1}{B} \sum_{b=1}^B \hat{\theta}_n^{(b)}$   
**Cramer Rao LB**  
 $\text{Var}(\hat{\theta}) \geq \frac{1}{E[\frac{\partial^2 \log p(x|\theta)}{\partial \theta^2}]}$   
 $\hat{\theta}$  is unbiased  
**Bayes Optimal Classifier** For 0-1 loss  
 $\hat{c}(x) = \arg \max_c P(c | x)$   
 $\hat{c}(x) = \arg \max_c \int P(y | x) \mathbb{I}(y \neq c) dy$   
 $\hat{c}(x) = \arg \max_c \int P(y | x) \mathbb{I}(y \neq c) dy$   
**SVD For linear regression**  
 $X^T X = U D V^T$   $U \in \mathbb{R}^{n \times n}$   $V \in \mathbb{R}^{d \times d}$   
 $X^T y = U D V^T U^T y = U D V^T U^T y$   
 $U^T y = U^T U D V^T U^T y = U^T y$   
 $V V^T = I_d$

**Classification**  
 GP:  $a(k+1) = a(k) - \eta \nabla_{a(k)}$   
 SGD  
 Newton optimiser  
 $a(k+1) = a(k) - H^{-1} \nabla J$   
 Linear sign( $a^T x$ ) for binary classification  
**Perceptron** Loss:  $J(a) = \sum_i \max(0, a^T x_i)$   
 $\rightarrow a(k+1) = a(k) + \eta (y_i - a(k)^T x_i) x_i$   
 Converges if data is separable.  
**Fisher's Linear Discriminant Analysis**  
 Maximize dist. of the means of the projected classes to find projection plane separating them best.  
 Projected mean for class  $a$   $\hat{\mu}_a = \frac{1}{n_a} \sum_{i \in a} w^T x_i$   
 Dist. of projected means for 2 classes  
 $J(w) = \frac{1}{2} (w^T \mu_1 - w^T \mu_2)^2$   
 Fisher's criterion:  $J(w) = \frac{(w^T \mu_1 - w^T \mu_2)^2}{w^T \Sigma w}$   
**Nearest neighbour**  
 LOA for multiple-class: One vs-rest  
**SVM:**  
 Primal constrained:  
 $\min_w w^T w + C \sum_i \xi_i$  s.t.  $y_i w^T x_i \geq 1 - \xi_i$   
 Primal unconstrained  
 $\min_w w^T w + C \sum_i \xi_i$  s.t.  $y_i w^T x_i \geq 1 - \xi_i$   
 Dual  
 $\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j x_i^T x_j$  s.t.  $0 \leq \alpha_i \leq C$   
 Lagrangian multiplier  
 $L(w, \mu, \alpha) = f(w) + \sum_i \mu_i (g_i(w) - 1) + \sum_i \alpha_i (y_i w^T x_i - 1 + \xi_i)$   
 General  
 $\min_{w, \mu, \alpha} L(w, \mu, \alpha) = f(w) + \sum_i \mu_i (g_i(w) - 1) + \sum_i \alpha_i (y_i w^T x_i - 1 + \xi_i)$   
 Strong dual  
**Kernel SVM**  
 $\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(x_i, x_j)$  s.t.  $0 \leq \alpha_i \leq C$   
 $y = \text{sign}(\sum_i \alpha_i y_i k(x, x_i))$   
**Logistic**  
 $P(y = 1 | x, w) = \frac{\exp(w^T x)}{1 + \exp(w^T x)}$   
**Multiclass Softmax**  
 $P(y = i | x, w) = \frac{\exp(w_i^T x)}{\sum_j \exp(w_j^T x)}$

**Multiclass SVM** still a hard margin.  
 $\min_{w, \gamma} \frac{1}{2} w^T w + C \sum_i \xi_i$  s.t.  $\forall y_i \in Y$   
 $(w^T y_i) - w^T x_i \geq 1 - \xi_i$   
**Structural SVM:**  
 $\min_{w, \gamma} \frac{1}{2} \|w\|_2^2 + \frac{1}{n} \sum_i \xi_i$  s.t.  $y_i \in Y$   
 where  $H(x) = \max_{y \in Y} (y^T x - \phi(x, y))$   
**XOR**  
 $K(x, y) = (1 + x_1 y_1 + x_2 y_2)^2$   
 $\frac{\partial L}{\partial w} = w - \sum_i \alpha_i y_i x_i$   
 $\frac{\partial L}{\partial \xi_i} = 0$   
 $\frac{\partial L}{\partial \xi_i} = 0$   
 $\frac{\partial L}{\partial \xi_i} = 0$   
 $L(x) = \frac{1}{2} w^T w + C \sum_i \xi_i$   
 $= \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j x_i^T x_j$   
 $= \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j x_i^T x_j$   
 Dual max  $V(\alpha) =$   
 subject to  $C \geq \alpha_i > 0$   
 $\alpha_i \leq C$   
**Primal for SSVM**  
 $\min_w \frac{1}{2} w^T w + C \sum_i \xi_i$   
 s.t.  $w^T \psi(z_i, y_i) - w^T \psi(z_i, y_i) \geq 1 - \xi_i$   
**Logistic:**  
 $L(x) = \frac{1}{2} w^T w - \sum_i \alpha_i [z_i^T (w y_i + w_0) - 1] - \sum_i \xi_i$   
 $\frac{\partial L}{\partial w} = 0$   
 $\frac{\partial L}{\partial w} = 0$   
 Dual  $\leq$  problem:  
 $\max_{\alpha} V(\alpha) = \sum_i \alpha_i - \sum_{i,j} \alpha_i \alpha_j z_i^T z_j$   
 s.t.  $\forall_i \alpha_i \geq 0$  and  $\sum_i \alpha_i y_i = 0$   
**Soft margin**  
 $w^* = \sum_i \alpha_i z_i y_i$   
 $w_0^* = -\frac{1}{2} (\min_i w_i^T y_i + \max_i w_i^T y_i)$   
**Primal:**  
 $\min_w \frac{1}{2} w^T w + C \sum_i \xi_i$   
 s.t.  $z_i^T (w y_i + w_0) \geq 1 - \xi_i$   
**Log:**  
 $L(\alpha, \beta) = \frac{1}{2} w^T w + C \sum_i \xi_i - \sum_i \alpha_i [z_i^T (w y_i + w_0) - 1] - \sum_i \xi_i$   
 $\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow C - \alpha_i - \beta_i = 0$   
 Primal is the same for this. But  $C \geq \alpha_i \geq 0$   
**Multiclass SVM w/ classes**  
 $\min_{w, \gamma} \frac{1}{2} w^T w + C \sum_i \xi_i$  s.t.  $\forall y_i \in Y$



