

- Validation
 - training set
 - test set
 - outliers, overfitting
 - want generalise
 - hyperparameter - control fit
 - time w/ validation set
- Supervised learning - classification, regression
- unsupervised learning - clustering, dimensionality reduction
- Classification - sample of observations w/ features
 - nearest neighbors
 - k-nearest neighbors
 - decision boundary
 - decision/predictor/discriminant function
 - $f(x) > 0$ SEC, $f(x) < 0$ N class
 - linear classifier - line/plane/hyperplane
 - inner product \Rightarrow dot product: linear!
 - Euclidean norm $\|x\| = \sqrt{x \cdot x}$
 - $\cos \theta = \frac{x \cdot y}{\|x\| \|y\|}$
 - $H = \{x: w \cdot x = a\}$ hyperplane $f(x) = w \cdot x + b$
 - $w \cdot x \geq a$ to all H vectors
 - $w \cdot x + a \Rightarrow$ signed distance: $\|w\|^{-1}$
 - linearly separable \Rightarrow plane exists
- Centroid Method
 - MC of all clusters in C , M_C not in C
 - $f(x) = (x - M_C) \cdot x - (M_C - M_C) \cdot \frac{M_C M_C}{2}$
 - good for Gaussian
- Perceptron Algorithm
 - slow, but correct for lin sep
 - use gradient descent
 - 1.0 C, -1.0 N
 - only origin H
 - find $w \Rightarrow x_1 \cdot w \geq 0$ if in C
 - aka $y_i \cdot x_i \cdot w \geq 0$ constraint
 - loss func $L(z, y) = \begin{cases} 0 & \text{if } z \geq y \\ -y + z & \text{else} \end{cases}$
 - risk func $R(w) = \frac{1}{n} \sum L(x_i \cdot w, y_i)$
 - optimization prob - find w min $R(w)$
 - x-space $w \cdot x \geq 0$, x
 - w-space $w \cdot x \geq 0$
 - $\nabla R(w) = - \sum y_i x_i$
 - $w \leftarrow w + \epsilon \nabla R(w)$ gradient descent $O(n)$
 - stochastic - pick one
 - add fictitious dimension: $f(x) = \begin{pmatrix} x \\ 1 \end{pmatrix}$
 - online alg - can learn while running
 - Convergence Theorem - $O(n^2)$ if ρ small $\|w\|$ margin
- Maximum Margin Classifier (SVM) (hard margin)
 - margin - distance H to pt closest
 - $y_i (w \cdot x_i + b) \geq 1 \Rightarrow$ margin: $\frac{1}{\|w\|} \cdot x_i \cdot \frac{1}{\|w\|} \geq \frac{1}{\|w\|}$
 - Find w, b , minimize $\|w\|^2$
 - quadratic prog. dsl dim. n constraints
- Soft Margin SVM
 - hard SVM - fail if not lin sep, but w/ outliers
 - add slack variables $\xi_i \geq 0$
 - $y_i (w \cdot x_i + b) \geq 1 - \xi_i$
 - margin = $\frac{1}{\|w\|}$
 - add bias term
 - min $\|w\|^2 + C \sum \xi_i$
 - small C - max $\|w\|$, underfit, linear, max fit
 - big C - imp slack, overfit, bad soln. $\xi_i = 0$
 - nonlinear?
 - make nonlinear features! hyperplane?
 - ex. add $\|x\|^2$
 - quadratic - histo dim. emv. solution
 - slow? kernel trick
 - edge detect features method?
- whitening data - mean=0, var=1
 - subtract mean $\mu(0, I)$
 - divide by std

- Optimization Alg
 - unconstrained - obj fun $f(w)$
 - smooth if $\nabla f(w)$ cont for
 - global/local min/max
 - convex \Rightarrow no above line
 - smooth fct
 - gradient descent (convergence)
 - newton's method
 - multi-step conjugate gradient
 - matmult - BFGS
 - line search - optimize in 1D
 - constrained - smooth
 - subject to $g(w) = 0$
 - LP - $w \in \text{min } c \cdot w$
 - subj to $Aw \leq b$
 - feasible region - convex polytope
 - optimum - extreme point
 - alg - simplex, interior point
 - QP $f(w) = w^T Q w + c^T w$
 - $Aw \leq b$
 - Q is PSD (convex)
 - alg - Simplex, SMO, coordinate descent
- Decision Theory \Rightarrow Risk min
 - X sample, Y class
 - $P(X|Y) = \frac{P(X, Y)}{P(Y)}$ joint w/ multivar
 - $P(a|b, c) = \frac{P(a, b, c)}{P(b, c)}$
 - posterior
 - loss function compare false res, pos diff (asym)
 - 0-1 loss func is symmetric
 - let r be decision, def risk
 - $R(r) = E[L(r(X), Y)] = \sum_x L(r, x) P(X=x)$
 - Bayes decision rule/classifier min $R(r)$ w/ Bayes rule
 - basically, pick class largest posterior prob
 - PDF $f(x) = \int_{-\infty}^{\infty} f(x, y) dy = 1$
 - $E(X) = \int_{-\infty}^{\infty} x f(x) dx$
 - $\sigma^2 = E[(X - \mu)^2] = E(X^2) - \mu^2$
 - Gaussian Discriminant Analysis
 - all class form Gaussian $X \sim N(\mu, \Sigma)$: $f(x) = \frac{1}{(2\pi)^d} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$
 - Bayes' rule $Q_C(x) = \ln \left(\frac{1}{(2\pi)^d} f_C(x) \pi_C \right)$
 - $\lim_{\Sigma \rightarrow 0} \text{PDF} \Rightarrow$ point
 - $\lim_{\Sigma \rightarrow \infty} \text{PDF} \Rightarrow$ uniform
 - $\frac{\|x - \mu\|^2}{2\sigma^2}$
 - $-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)$
 - precision matrix
 - covariance
 - MLE
 - pick model, param that best fit p max $L(p)$
 - $\frac{dL}{dp} = 0$, check $\frac{d^2L}{dp^2} < 0$
 - Gaussian $L(\cdot) = f_C(x) f_N(x) \dots f_{K-1}(x)$ $\hat{\mu} = \frac{1}{n} \sum x_i$
 - can also max log likelihood $\hat{\sigma}^2 = \frac{1}{n} \sum \|x_i - \hat{\mu}\|^2$
 - Eigenvectors
 - $A v = \lambda v$, just scalar, $A^k v = \lambda^k v$
 - Spectral decomp real sym $A = V \Lambda V^T$ "eigenval" each λ^i
 - general $A = a b a^T$
 - symm PD $x^T A x > 0$
 - PSD $x^T A x \geq 0$
 - QDA
 - $\hat{\Sigma}_C = \frac{1}{n_C} \sum (x_i - \hat{\mu}_C)(x_i - \hat{\mu}_C)^T$ empirical cov
 - min quadratic decision func
 - $Q_C(x) = -\frac{1}{2}(x - \hat{\mu}_C)^T \hat{\Sigma}_C^{-1}(x - \hat{\mu}_C) - \frac{1}{2} \ln |\hat{\Sigma}_C| + \ln \pi_C$
 - LDA
 - one $\hat{\Sigma} = \frac{1}{n} \sum \sum (x_i - \hat{\mu}_C)(x_i - \hat{\mu}_C)^T$ for all
 - $Q_C(x) = \underbrace{M_C^T \sum^{-1} x}_{w \cdot x} - \frac{1}{2} M_C^T \sum^{-1} M_C + \ln \pi_C$
- not perfect bc still estimating, data w/ perfect Gaussian
 - X rank $n \times d$ \Rightarrow design matrix
 - sphering $w = \frac{1}{\sqrt{\lambda}} \text{Var}(R)^{-1/2}$
 - $\lambda = \frac{1}{\sigma^2}$

- Regression
 - regression fun h
 - (1) $w \cdot x + b$ - linear
 - (2) polynomial
 - (3) logistic $s(w \cdot x + b)$ $s(z) = \frac{1}{1 + e^{-z}}$
 - loss fun L
 - (A) $(z - y)^2$ - squared
 - (B) $|z - y|$ - abs
 - (C) $-y \ln z - (1 - y) \ln(1 - z)$ - log-likelihood / cross entropy
 - cost fun J
 - (a) $\frac{1}{n} \sum L(h(x_i), y_i)$ - mean
 - (b) $\max_{\theta} \sum L(h(x_i), y_i)$ - max
 - (c) $\sum_{i=1}^n L(h(x_i), y_i)$ - weighted sum
 - (d) $\frac{1}{n} \sum L(h(x_i), y_i) + \lambda \|w\|^2$ L_2 penalized / regularized
 - (e) $\frac{1}{n} \sum L(h(x_i), y_i) + \lambda \|w\|_1$ L_1
 - least squares linear reg $Y = A \cdot x + a$
 - min $\|Xw - y\|^2 = \text{RSS}(w)$
 - $w = (X^T X)^{-1} X^T y$
 - full β unconstrained, positive values
 - logistic reg $\exists t \in \mathbb{R}$
 - $\theta w^T = -X^T (y - s(Xw))$
 - Newton's Method - GD but smarter step
 - $\nabla^2 J(w) = \text{Hessian}$
 - $w \leftarrow w + (-\nabla^2 J(w))^{-1} \nabla J(w)$
 - much gut stab
 - ROC curve (Receiver operating characteristic)
 - rate false pos vs true pos
 - $\text{IROC} = \int_0^1 \text{ROC}(t) dt$
 - $\text{IROC} = 0.5$ if random
 - $\text{IROC} = 1$ if perfect
 - $\text{IROC} = 0.5$ if random
 - $\text{IROC} = 1$ if perfect
 - empirical distribution - discrete $P(X=x)$
 - empirical risk - $\hat{R}(w) = \frac{1}{n} \sum L(h(x_i), y_i)$
 - Bias Variance Decomposition
 - bias - error from bad model?
 - variance - error from noise
 - $E(h(x)) - g(x)$ bias
 - $\text{Var}(h(x))$ variance
 - $\text{Var}(y)$ irreducible error
 - Ridge Regression (Tikhonov)
 - locally line of optimal, encourage small $\|w\|$
 - unique sol. PD
 - $I + \lambda I$
 - $X^T X + \lambda I \Rightarrow$ PD!
 - balances bias & var
 - any design prob, just whether data
 - MAP (maximum a posteriori)
 - LASSO
 - $1 + \lambda \cdot \ell$ - max sparse
 - $\nabla_x (y - z) = (y - z) z + (y - z)^2$
 - $\nabla_x x y = (y - z) y^T + x \nabla_x y$
 - $\nabla_x f(y) = (y - z) y^T (y - z)$
 - $\nabla_x (y - z) = (y - z) y^T (y - z)$
 - $\nabla_x A x = A^T$ $\nabla_x x = I$

