

# Machine Learning Cheat Sheet

## General

### Dataset and Definitions

$\mathcal{D}$  is a set of training examples, the  $n$ -th Training Example ( $n = 1, 2, \dots, N$ ), of this set is:

$$\mathbf{x}_n = [x_{n1} \quad x_{n2} \quad \dots \quad x_{nD}]$$

We write all  $N$  training examples in a matrix:

$$\mathbf{X} = [\mathbf{x}_1; \quad \mathbf{x}_2; \quad \dots; \quad \mathbf{x}_N]$$

In supervised learning, you are also given a set of observations corresponding to the training examples:

$$\mathbf{y} = [y_1 \quad y_2 \quad \dots \quad y_N]^T$$

We assume a *true* model behind the data, the observations are noisy versions of this ground truth:

$$y = y_{true} + \epsilon.$$

The final goal is to predict a  $\hat{y} = f(\hat{\mathbf{x}})$  given any  $\hat{\mathbf{x}}$ .

### Distributions

Multivariate gaussian distribution:  $\mathcal{N}(\mathbf{X}|\mu, \Sigma)$

$$\implies p(\mathbf{X} = \mathbf{x}) =$$

$$(2\pi)^{-d/2} |\Sigma|^{-1/2} \exp[-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)]$$

Gaussian distribution:  $\mathcal{N}(X|\mu, \sigma^2)$

$$\implies p(X = x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2}(\frac{x-\mu}{\sigma})^2)$$

Poisson distribution:  $\mathcal{P}(X|\lambda)$

$$\implies p(X = k) = \frac{\lambda^k}{k!} \exp(-\lambda)$$

### Properties

**Bayes rule:**  $p(A, B) = p(A|B)p(B) = p(B|A)p(A)$   
A matrix  $\mathbf{M}$  is **positive semidefinite**  $\iff \forall$  nonzero vector  $\mathbf{a}$ , we have  $\mathbf{a}^T \mathbf{M} \mathbf{a} \geq 0$

If  $\mathbf{V}$  symmetric positive definite, then for all  $\mathbf{P} \neq 0$ ,  $\mathbf{P}^T \mathbf{V} \mathbf{P}$  is positive semi-definite (and even positive definite if  $\mathbf{P}$  is not singular).

**Jensen's inequality** applied to log:

$$\log(\mathbb{E}[X]) \geq \mathbb{E}[\log(X)]$$

$$\implies \log(\sum_x x \cdot p(x)) \geq \sum_x p(x) \log(x)$$

**Matrix inversion lemma:**

$$(\mathbf{PQ} + \mathbf{I}_N)^{-1} \mathbf{P} = \mathbf{P}(\mathbf{QP} + \mathbf{I}_M)^{-1}$$

Useful derivative:  $\frac{\partial \mathbf{x}^T \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{B} + \mathbf{B}^T) \mathbf{x}$

Marginal and Conditional Gaussians:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mu, \Lambda^{-1})$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$$

$$\Downarrow$$

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mu + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\Lambda^{-1}\mathbf{A}^T)$$

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\Sigma\{\mathbf{A}^T \mathbf{L}(\mathbf{y} - \mathbf{b}) + \Lambda\mu\}, \Sigma)$$

$$\text{where } \Sigma = (\Lambda + \mathbf{A}^T \mathbf{L} \mathbf{A})^{-1}$$

## Optimization methods

### Grid Search

Simply try values for all parameters at regular intervals. Complexity:  $\mathcal{O}(M^D ND)$ , where  $M$  is the number of values tried in each dimension.

## Gradient Descent

Update rule:  $\beta^{(k+1)} = \beta^{(k)} - \alpha \frac{\partial \mathcal{L}(\beta^{(k)})}{\partial \beta}$

Complexity:  $\mathcal{O}(IND)$  where  $I$  is the number of iterations we take.

$\alpha$  is a parameter that needs to be chosen carefully.

The gradient for MSE comes out as:

$$\frac{\partial \mathcal{L}}{\partial \beta} = -\frac{1}{N} \tilde{\mathbf{X}}^T (\mathbf{y} - \tilde{\mathbf{X}}\beta)$$

### Newton's method

It uses second order derivatives information to converge faster (we approximate the objective function by a quadratic function rather than a line).

General rule:  $\beta^{(k+1)} = \beta^{(k)} - \alpha \mathbf{H}_k^{-1} \frac{\partial \mathcal{L}(\beta^{(k)})}{\partial \beta}$

where  $\mathbf{H}_k$  is the  $D \times D$  Hessian at step  $k$ :

$$\mathbf{H}_k = \frac{\partial^2 \mathcal{L}(\beta^{(k)})}{\partial \beta^2}$$

Newton's method is equivalent to solving many least-squares problems.

Complexity:  $\mathcal{O}(IND^2 + ID^3)$ , with the  $D^3$  cost coming from the inversion of  $\mathbf{H}_k$ .

## Expectation-Maximization

EM is a family of methods to find a MLE/MAP estimator if some of the data is missing (e.g. latent variables):

Given is some data  $\mathbf{x}$  and a model with some latent variables  $\mathbf{z}$  and a likelihood  $\mathcal{L}(\theta; \mathbf{x}, \mathbf{z})$  (e.g. a gaussian where we need to find the mean and the covariance). Goal is to find  $\beta = \arg \max_{\beta} \mathcal{L}_{lik}(\theta; x)$ , i.e. the MLE given only the data, without knowing the latent variables.

The problem is, that this likelihood can oftentimes not be optimized in closed form with respect to  $\theta$ . The solution is EM: Iteratively find better  $\theta$ , start with  $\theta_0$  and iterate with variable  $t$ :

- E-step calculates the  $\mathbb{E}$  of log likelihood, with respect to  $\mathbf{Z}$  given  $\mathbf{X}$  under the current estimate  $\theta_t$ :  $Q(\theta, \theta_t) = \mathbb{E}_{\mathbf{Z}|\mathbf{X}, \theta_t}[\log(p(\mathbf{X}, \mathbf{Z}|\theta)|X = x)]$

- M-step finds next estimate that maximizes:  $\theta_{t+1} = \arg \max_{\theta} Q(\theta, \theta_t)$

- Iterate until convergence.

## Regression

Simple linear regression:  $y_n \approx \beta_0 + \beta_1 x_{n1}$

Multiple linear regression:

$$y_n \approx f(\mathbf{x}_n) := \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_D x_{nD}$$

### Linear basis function model

We can create more complex models while staying in the linear framework by transforming the inputs  $X$  of dimensionality  $D$  through a function  $\phi: D \rightarrow M$ .

$$y_n = \beta_0 + \sum_{i=1}^M \beta_i \phi_i(\mathbf{x}_n) = \tilde{\phi}^T(\mathbf{x}_n^T) \beta.$$

The optimal  $\beta$  can be computed in closed form by

$$\beta = (\tilde{\Phi}^T \tilde{\Phi})^{-1} \tilde{\Phi}^T \mathbf{y} \text{ where } \tilde{\Phi} \text{ is a matrix with } N \text{ rows and the } n\text{-th row is } [1, \phi_1(\mathbf{x}_n), \dots, \phi_M(\mathbf{x}_n)].$$

But note this requires  $\tilde{\Phi}^T \tilde{\Phi}$  to be invertible (well-conditioned:  $\tilde{\Phi}$  full column-rank).

$$\text{Ridge regression: } \beta_{\text{ridge}} = (\tilde{\Phi}^T \tilde{\Phi} + \lambda \mathbf{I})^{-1} \tilde{\Phi}^T \mathbf{y}$$

## Cost functions

Cost function / Loss:  $\mathcal{L}(\beta) = \mathcal{L}(\mathcal{D}, \beta)$

Mean square error (MSE):  $\frac{1}{2N} \sum_{n=1}^N [y_n - f(\mathbf{x}_n)]^2$

MSE matrix formulation:  $\frac{1}{2N} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)$

Mean absolute error (MAE):  $\frac{1}{2N} \sum_{n=1}^N |y_n - f(\mathbf{x}_n)|$

Huber loss:  $\mathcal{L}_{\delta}(a) = \begin{cases} \frac{1}{2}a^2 & \text{for } |a| \leq \delta, \\ \delta(|a| - \frac{1}{2}\delta), & \text{otherwise.} \end{cases}$

Root mean square error (RMSE):  $\sqrt{2 * \text{MSE}}$

Epsilon insensitive ("hinge loss", used for SVM):

$$\mathcal{L}_{\epsilon}(y, \hat{y}) = \begin{cases} 0 & \text{if } |y - \hat{y}| \leq \epsilon, \\ |y - \hat{y}| - \epsilon, & \text{otherwise.} \end{cases}$$

## Least squares

Complexity:  $\mathcal{O}(ND^2 + D^3)$  The gradient of the MSE set to zero is called the normal equation:

$$\mathbf{X}^T \mathbf{X} \beta - \mathbf{X}^T \mathbf{y} = 0$$

We can solve this for  $\beta$  directly (by matrix manipulation)  $\beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

## Classification

Logistic Function:  $\sigma(t) = \frac{\exp(t)}{1 + \exp(t)}$

Derivative:  $\frac{\partial \sigma(t)}{\partial t} = \sigma(t)[1 - \sigma(t)]$

Classification with linear regression: Use  $y = 0$  as class  $\mathcal{C}_1$  and  $y = 1$  as class  $\mathcal{C}_2$  and then decide a newly estimated  $y$  belongs to  $\mathcal{C}_1$  if  $y < 0.5$ .

## Logistic Regression

Formally, the outcomes  $Y_n$  are described as being Bernoulli-distributed data, where each outcome is determined by an unobserved probability  $p_i$  that is specific to the outcome at hand, but related to the explanatory variables:  $\mathbb{E}[Y_n | x_{1,i}, \dots, x_{m,i}] = p_i$ . The model can then be formulated as:

$$p_i = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_{1,i} + \dots + \beta_k x_{k,i})}}$$

There's no closed form, we can use gradient descent with gradient  $\frac{\partial \mathcal{L}(\beta)}{\partial \beta} = \tilde{\mathbf{X}}^T [\sigma(\tilde{\mathbf{X}}\beta) - \mathbf{y}]$

### Assumptions of Logistic Regression

- Target variable is binary.
- Observations are independent.
- Predictors not highly correlated.
- Predictors and log odds are linearly related.
- Large sample size ('Rule of 10'), i.e. minimum 10 cases with least frequent outcome for each predictor in the model.
- Not assuming:* linearity, normality, homoscedasticity and measurement level.

## Generalized linear model

The GLM consists of three elements:

- A probability distribution from the exponential family.
- A linear predictor  $\hat{y} = \mathbf{X}\beta$ .
- A link function  $g$  such that  $E(y) = \mu = g^{-1}(\eta)$ .

In a generalized linear model (GLM), each outcome of the dependent variables  $y$  is assumed to be generated from a particular distribution in the exponential family, a large range of probability distributions that includes the normal, binomial, Poisson and gamma distributions, among others.

## Cost functions

$$\text{RMSE: } \sqrt{\frac{1}{N} \sum_{n=1}^N [y_n - \hat{p}_n]^2}$$

$$0\text{-}1 \text{ Loss: } \frac{1}{N} \sum_{n=1}^N \delta(y_n, \hat{y}_n)$$

Log-Loss:

$$-\frac{1}{N} \sum_{n=1}^N y_n \log(\hat{p}_n) + (1 - y_n) \log(1 - \hat{p}_n)$$

## Probabilistic framework

The Likelihood Function maps the model parameters to the probability distribution of  $\mathbf{y}$ :

$\mathcal{L}_{lik}$ : parameter space  $\rightarrow [0; 1]$   $\beta \mapsto p(\mathbf{y} | \beta)$  An underlying  $p$  is assumed before. If the observed  $y$  are IID,  $p(\mathbf{y} | \beta) = \prod_n p(y_n | \beta)$ .

$\mathcal{L}_{lik}$  can be viewed as just another fitness function. Maximum likelihood then simply chooses the parameters  $\beta$  such that observed data is most likely.  $\beta = \arg \max_{\beta} L(\beta)$

Assuming different  $p$  is basically what makes this so flexible. We can chose e.g.:

Gaussian $p$	$\mathcal{L}_{lik} \hat{=}$	$-\mathcal{L}_{MSE}$
Laplace $p$	$\mathcal{L}_{lik} \hat{=}$	$-\mathcal{L}_{MAE}$

It is a sample approximation of the expected likelihood:  $\mathcal{L}_{lik}(\beta) \approx E_y[p(y | \beta)]$

## Bayesian methods

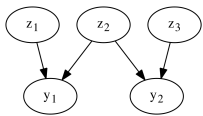
Bayes rule:  $p(A, B) = p(A|B)p(B) = p(B|A)p(A)$

The **prior**  $p(\mathbf{f}|\mathbf{X})$  encodes our prior belief about the "true" model  $\mathbf{f}$ . The **likelihood**  $p(\mathbf{y}|\mathbf{f})$  measures the probability of our (possibly noisy) observations given the prior.

Least-squares tries to find model parameters  $\beta$  which maximize the likelihood. Ridge regression maximizes the **posterior**  $p(\beta|\mathbf{y})$

Bayesian networks

We can use a Directed Acyclic Graph (DAG) to define a joint distribution of events. For example, we can express the relationship between *latent factors* (possible “causes”)  $z_i$  and *observations* (results)  $y_i$ :



This example can be factorized as follows:

$$p(y_1, y_2, z_1, z_2, z_3) = p(y_1|z_1, z_2)p(y_2|z_2, z_3)p(z_1)p(z_2)p(z_3)$$
  
We can then obtain the distribution over latent factors ( $z_i$ ) by marginalizing over the unknown variables:  $p(z_1, z_2, z_3|y_1, y_2) = \frac{\text{joint}}{p(y_1, y_2)}$   
 $\implies p(z_1|y_1, y_2) = \sum_{z_2, z_3} \frac{\text{joint}}{p(y_1, y_2)}$

Belief propagation

Belief propagation is a message-passing based algorithm used to compute desired marginals (e.g.  $p(z_1|y_1, y_2)$ ) efficiently. It leverages the factorized expression of the joint. Messages passed:

$$m_{z_i \rightarrow y_j} = p(z_i) \prod (\text{messages received except from } y_j)$$
  
$$m_{y_j \rightarrow z_i} = \sum_{z_k \neq z_i} p(y_j|z_k) \prod (\text{messages received except from } z_i)$$

Kernel methods

Kernels are another way to encode our prior belief about the relation of outputs to inputs. A kernel defines a distance measure, or “similarity” of two vectors. We define:

$(\mathbf{K})_{i,j} = \kappa(\mathbf{x}_i, \mathbf{x}_j) = \vec{\phi}(\mathbf{x}_i)^T \vec{\phi}(\mathbf{x}_j)$ .  
The  $\phi$  are not that important in the end, because we only use the Kernel as is. Sometimes it’s even impossible to write them down explicitly.

Linear	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
Polynomial	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + c)^d$
RBF	$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\ \mathbf{x}_i - \mathbf{x}_j\ ^2}{2\sigma^2}\right)$

Properties of a Kernel:

- $\mathbf{K}$  should be symmetric:  $\mathbf{K}^T = \mathbf{K}$
- $\mathbf{K}$  should be positive semi-definite:  $\forall$  nonzero vector  $\mathbf{a}, \mathbf{a}^T \mathbf{K} \mathbf{a} \geq 0$ .

Gaussian Process

The predicting function  $f$  is interpreted as a random variable with jointly gaussian prior:  $\mathcal{N}(f|\mathbf{0}, \mathbf{K})$ .  
Defining the Kernel matrix  $\mathbf{K} = \kappa(\mathbf{X}, \mathbf{X})$  defines the prior. The key idea is, that if  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are deemed by the kernel to be similar, then we expect the output of  $f$  at those points to be similar, too.

We can sample functions from this random variable  $f$  and we can use prior + measurements to generate predictions.

If we have measurements  $\mathbf{y}$  available, we get a joint distribution with the  $\hat{\mathbf{y}}$  to be predicted:

$$\begin{bmatrix} \mathbf{y} \\ \hat{\mathbf{y}} \end{bmatrix} = \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} \kappa(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} & \kappa(\mathbf{X}, \hat{\mathbf{X}}) \\ \kappa(\hat{\mathbf{X}}, \mathbf{X}) & \kappa(\hat{\mathbf{X}}, \hat{\mathbf{X}}) \end{bmatrix}\right)$$
  
This can be conditioned on  $\mathbf{y}$  to find the PDF of  $\hat{\mathbf{y}}$ .  
Advantage: We output our prediction as probabilities (which represent uncertainty).

Neural Networks

A feed forward Neural Network is organized in  $K$  layers, each layer with  $M^{(k)}$  hidden units  $z_i^{(k)}$ .  
Activations  $a_i^{(k)}$  are computed as the linear combination of the previous layer’s terms, with weights  $\beta^{(k)}$  (one  $M^{(k-1)} \times 1$  vector of weights for each of the  $M^{(k)}$  activations). Activations are then passed through a (possibly nonlinear) function  $h$  to compute the hidden unit  $z_i^{(k)}$ .

$$\mathbf{x}_n \xrightarrow{\beta^{(1)}} a_i^{(1)} \xrightarrow{h} z_i^{(1)} \xrightarrow{\beta^{(2)}} \dots \mathbf{z}^{(K)} = \mathbf{y}_n$$

Backpropagation

It’s an algorithm which computes the gradient of the cost  $\mathcal{L}$  w.r.t. the parameters  $\beta^{(k)}$ .  
Forward pass: compute  $a_i, z_i$  and  $\mathbf{y}_n$  from  $\mathbf{x}_n$ .  
Backward pass: work out derivatives from outputs to the target  $\beta_i^{(k)}$ . Using the chain rule:

$$\delta^{(k-1)} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{(k-1)}} = \text{diag}[h'(\mathbf{a}^{(k-1)})] \mathbf{B}^{(k)T} \delta^{(k)}$$
  
$$\frac{\partial \mathcal{L}}{\partial \mathbf{B}^{(1)}} = \delta^{(1)} \mathbf{x}^T$$
  
$$\frac{\partial \mathcal{L}}{\partial \mathbf{B}^{(k)}} = \delta^{(k)} \mathbf{z}^{(k)T}$$

Regularization

NN are not *identifiable* (existence of many local optima), therefore the maximum likelihood estimator is not *consistent*.  
NN are universal density estimators, and thus prone to severe overfitting. Techniques used to reduce overfitting include early stopping (stop optimizing when test error starts increasing) and “weight decay” (i.e.  $L_2$  regularization).

Support Vector Machines

Search for the hyperplane separating the data such that the gap (margin) is biggest. It minimizes the following cost function (“hinge loss”):  
$$\mathcal{L}_{SVM}(\beta) = \sum_{n=1}^N [1 - y_n \phi_n \beta]_+ + \frac{\lambda}{2} \sum_{j=1}^M \beta_j^2$$
  
with  $[t]_+ = \max(0, t)$   
This is convex but not differentiable. In the dual, the same problem can be formulated as:  
$$\max_{\alpha \in [0; C]^N} \alpha^T \mathbf{1} - \frac{1}{2} \alpha^T \mathbf{Y} \mathbf{K} \mathbf{Y} \alpha, \quad \alpha^T \mathbf{y} = 0$$

Gaussian Mixture Models

In mixture models, the data is generated by a sum (a mix) of  $K$  models. For GMM, these are gaussian:

$$p(\mathbf{x}_i|\theta) = \sum_{k=1}^K \pi_k p_k(\mathbf{x}_i|\theta) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \underbrace{\mu_k, \Sigma_k}_{\theta})$$

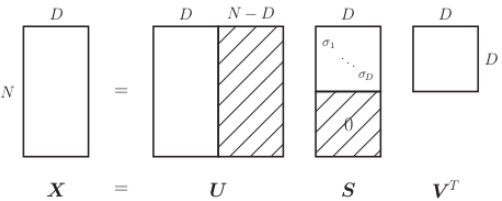
To use this for clustering, we first fit this mixture and then compute the posterior  $p(z_i = k|\mathbf{x}_i, \theta)$ . This yields *soft* cluster assignments.

PCA

Find the eigenvectors of the covariance matrix  $\mathbf{X}^T \mathbf{X}$  of the data. These form an orthonormal basis  $\{\mathbf{w}_1, \dots, \mathbf{w}_N\}$  for the data in the directions that have highest variance. One can then use the first  $L < D$  vectors to rebuild the data:  $\hat{\mathbf{x}}_i = \mathbf{W} \mathbf{z}_i = \mathbf{W} \mathbf{W}^T \mathbf{x}_i$ , with  $\mathbf{W} = [\mathbf{w}_1; \dots; \mathbf{w}_L]$ . This minimizes mean square error  $\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i - \hat{\mathbf{x}}_i^2$ .

SVD

The same as with PCA, we can do with SVD:



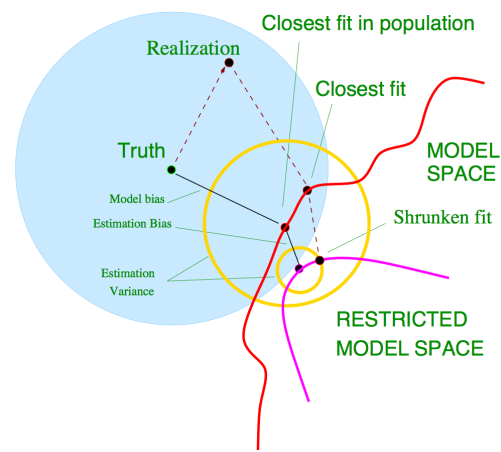
The singular vals of a  $N \times D$  matrix  $\mathbf{X}$  are the square roots of the eigenvalues of the  $D \times D$  matrix  $\mathbf{X}^T \mathbf{X}$

Concepts

Convexity

$f$  is called convex f:  $\forall x_1, x_2 \in X, \forall t \in [0, 1] : f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$ .  
Sum of two convex functions is convex. Composition of a convex function with a convex, nondecreasing function is convex. Linear, exponential and  $\log(\sum \exp)$  functions are convex.

Bias-Variance Decomposition



Bias-variance comes directly out of the test error:

$$\begin{aligned} \overline{teErr} &= E[(\text{observation} - \text{prediction})^2] = E[(y - \hat{y})^2] \\ &= E[(y - y_{true} + y_{true} - \hat{y})^2] \\ &= \underbrace{E[(y - y_{true})^2]}_{\text{var of measurement}} + E[(y_{true} - \hat{y})^2] \\ &= \sigma^2 + E[(y_{true} - E[\hat{y}] + E[\hat{y}] - \hat{y})^2] \\ &= \sigma^2 + \underbrace{E[(y_{true} - E[\hat{y}])^2]}_{\text{pred bias}^2} + \underbrace{E[(E[\hat{y}] - \hat{y})^2]}_{\text{pred variance}} \end{aligned}$$

	bias	variance
regularization	+	-
choose simpler model	+	-
more data	-	

Identifiability

We say that a statistical model  $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$  is identifiable if the mapping  $\theta \mapsto P_\theta$  is one-to-one:  
 $P_{\theta_1} = P_{\theta_2} \implies \theta_1 = \theta_2$  for all  $\theta_1, \theta_2 \in \Theta$ .

A non-identifiable model will typically have many local optima yielding the same cost, e.g.  
 $\mathcal{L}(W, Z) = \mathcal{L}(aW, \frac{1}{a}Z)$

Curse of dimensionality

With the dimensionality increase, every data point becomes arbitrarily far from every other data point and therefore the choice of nearest neighbor becomes random. In high dimension, data only covers a tiny fraction of the input space, making generalization extremely difficult. Or in other words, the volume of the space increases so fast that the available data become sparse.

## Primal vs. Dual

Instead of working in the **column space** of our data, we can work in the **row space**:

$$\hat{y} = \mathbf{X}\beta = \mathbf{X}\mathbf{X}^T\alpha = \mathbf{K}\alpha$$

where  $\beta \in \mathbb{R}^D$  and  $\alpha \in \mathbb{R}^N$  and (like magic)  $\mathbf{K}$  shows up, the Kernel Matrix.

Representer Theorem: For any  $\beta$  minimizing

$$\min_{\beta} \sum_{n=1}^N \mathcal{L}(y_n, \mathbf{x}_n^T \beta) + \sum_{d=1}^D \lambda \beta_d^2$$

there exists an  $\alpha$  such that  $\beta = \mathbf{X}^T \alpha$ .

When we have an explicit vector formulation of  $\beta$ , we can use the matrix inversion lemma to get to the dual. E.g. for ridge regression:

$$\beta = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{X}^T \underbrace{(\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I}_N)^{-1}}_{\alpha} \mathbf{y}$$

In optimization, we get to the dual like this:

$$\begin{array}{ccc} \min_{\beta} g(\beta) & \xrightarrow{(1)} & \min_{\beta} \max_{\alpha} G(\beta, \alpha) \\ & & \downarrow (2) \\ \max_{\alpha} g^*(\alpha) & \xleftarrow{(3)} & \max_{\alpha} \min_{\beta} G(\beta, \alpha) \\ & & g^*(\alpha) \end{array}$$

## Consistency

An estimator is said to be consistent, if it eventually recovers the true parameters that generated the data as the sample size goes to infinity. Of course, this only makes sense if the data actually comes from the specified model, which is not usually the case. Nevertheless, it can be a useful theoretical property.

## Efficiency

An estimator is called efficient if it achieves the Cramer-Rao lower bound:  $\text{Var}(\beta) \geq 1/I(\beta)$ , where  $I$  is the Fisher information.

## Occam's Razor

It states that among competing hypotheses, the one with the fewest assumptions should be selected. Other, more complicated solutions may ultimately prove correct, but in the absence of certainty, the fewer assumptions that are made, the better.

TODO: K-fold cross-validation, definition of Test-Error, Train-Error

TODO: statistical goodness (robust to outliers, ...) vs. computational goodness (convex, low computational complexity, ...) tradeoff. No free lunch theorem.

## Decision Trees and Random Forest

TODO: Decision Trees and Random Forests and Model averaging

## Recommender Systems

- Non-personalized
- Personalized

Technique	Background	Input	Process
Collaborative	Ratings from $U$ of items in $I$ .	Ratings from $u$ of items in $I$ .	Identify users in $U$ similar to $u$ , and extrapolate from their ratings of $i$ .
Content-based	Features of items in $I$ .	$u$ 's ratings of items in $I$ .	Generate a classifier that fits $u$ 's rating behavior and use it on $i$ .
Demographic	Demographic information about $U$ and their ratings of items in $I$ .	Demographic information about $u$ .	Identify users that are demographically similar to $u$ , and extrapolate from their ratings of $i$ .
Utility-based	Features of items in $I$ .	A utility function over items in $I$ that describes $u$ 's preferences.	Apply the function to the items and determine $i$ 's rank.
Knowledge-based	Features of items in $I$ . Knowledge of how these items meet a user's needs.	A description of $u$ 's needs or interests.	Infer a match between $i$ and $u$ 's need.

Technique	Pluses	Minuses
Collaborative filtering (CF)	A. Can identify cross-genre niches. B. Domain knowledge not needed. C. Adaptive: quality improves over time. D. Implicit feedback sufficient	I. New user ramp-up problem J. New item ramp-up problem K. "Gray sheep" problem L. Quality dependent on large historical data set. M. Stability vs. plasticity problem
Content-based (CN)	B, C, D	I, L, M
Demographic (DM)	A, B, C	I, K, L, M N. Must gather demographic information
Utility-based (UT)	E. No ramp-up required F. Sensitive to changes of preference G. Can include non-product features	O. User must input utility function P. Suggestion ability static (does not learn)
Knowledge-based (KB)	E, F, G H. Can map from user needs to products	P Q. Knowledge engineering required.

## ML-Architecture

1. Data Sourcing: Preferably into normalized, indexed database
2. Feature Engineering: SQL queries or transformations to build features
3. Feature Loader: Function with parameters to load different Features, can be part of pipeline
4. Imputer & Normalizer or other feature preprocessing
5. ML-Algorithm

## One-Hot Encoding

Is a binarization of categorical variables. Each category gets its own dimension. For example, go from this...

CompanyName	Categoricalvalue	Price
VW	1	20000
Acura	2	10011
Honda	3	50000
Honda	3	10000

to the following...

VW	Acura	Honda	Price
1	0	0	20000
0	1	0	10011
0	0	1	50000
0	0	1	10000

## Credits

Most material was taken from the lecture notes of Prof. Emtiyaz Khan.

Cost functions figure from Patrick Breheny's slides.

Biais-variance decomposition figure from Hastie,

Tibshirani, Friedman, *The Elements of Statistical*

*Learning*. The SVD figure from Kevin P. Murphy,

*Machine Learning, A Probabilistic Perspective*.

Rendered May 16, 2018. Written by Dennis Meier and Merlin Nimier-David.

© Dennis Meier. This work is licensed under the Creative Commons Attribution-ShareAlike 3.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-sa/3.0/> or send a letter to Creative Commons, 444 Castro Street, Suite 900, Mountain View, California, 94041, USA.

