

# PCML Cheat Sheet

## 1 Math Prerequisites

- Bayes rule

$$p(A, B) = \underbrace{p(A|B)}_{\text{Lik. Prior}} \underbrace{p(B)}_{\text{Post Marg. Lik.}}$$

- Gaussian distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{2\pi|\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right)$$

- Production of independent variables:

$$V(XY) = E(X^2)E(Y^2) - [E(X)]^2[E(Y)]^2$$

- Log-properties

$$\log(mn) = \log(m) + \log(n)$$

$$\log(m^n) = n \log(m)$$

- Covariance matrix of a data vector  $\mathbf{x}$

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - E(\mathbf{x}))(\mathbf{x}_n - E(\mathbf{x}))^T$$

## 1.1 Convexity

- A function is convex when a line joining two points never intersects with the function anywhere else.
- A function  $f(\mathbf{x})$  is convex, if for any  $\mathbf{x}_1, \mathbf{x}_2 \in \mathbf{X}$  and for any  $0 \leq \lambda \leq 1$ , we have :

$$f(\lambda \mathbf{x}_1 + (1-\lambda)\mathbf{x}_2) \leq \lambda f(\mathbf{x}_1) + (1-\lambda)f(\mathbf{x}_2)$$

- A function is strictly convex if the inequality is strict.
- A convex function has only one global minimum.
- Sums of convex functions are also convex.
- The Hessian is related to the convexity of a function: a twice differentiable function is convex if and only if the Hessian is positive definite.
- The Hessian of a convex function is positive semi-definite and for a strictly-convex function it is positive definite.
- The Hessian matrix of a function

$$\mathbf{H}_{i,j} = \frac{d^2 f}{dx_i dx_j}$$

## 1.2 Linear Algebra

- Column  $\mathbf{x} \in R^n$ , rows  $\mathbf{x}^T$ , matrix  $\mathbf{A} \in R^{m \times n}$
- $\mathbf{x}^T \mathbf{x}$  is a scalar,  $\mathbf{x} \mathbf{x}^T$  is a matrix
- $\mathbf{A}^{-1}$  exist if  $\mathbf{A}$  is full rank
- **Condition number** of a function measures how much the output value can change for a small change in the input argument. A matrix with a high condition number is said to be **ill-conditioned**. If  $\mathbf{A}$  is normal ( $A^T A = A A^T$ ) then

$$k(\mathbf{A}) = \left| \frac{\lambda_{\max}(\mathbf{A})}{\lambda_{\min}(\mathbf{A})} \right|$$

- A positive definite matrix is **symmetric** with all positive eigenvalues
- The real symmetric  $N \times N$  matrix  $\mathbf{V}$  is said to be **positive semidefinite** if

$$\mathbf{a}^T \mathbf{V} \mathbf{a} \geq 0$$

for any real  $N \times 1$  vector  $\mathbf{a}$ .

- The real symmetric  $N \times N$  matrix  $\mathbf{V}$  is said to be **positive definite** if

$$\mathbf{a}^T \mathbf{V} \mathbf{a} > 0$$

for any real  $N \times 1$  vector  $\mathbf{a}$ .

- Cost of matrix inversion:  $O(n^3) \rightarrow O(n^{2.372})$
- Cost of determinant computation using LU decomposition:  $O(n^3)$

## 2 Cost functions

- Cost functions are used to learn parameters that explain the data well.
- It is essential to make sure that a global minimum exist  $\rightarrow$  lower bounded

**Mean square error (MSE):**

$$MSE(\boldsymbol{\beta}) = \sum_{n=1}^N (y_n - f(\mathbf{x}_n))^2$$

- MSE is **convex** thus it has only one global minimum value.
- MSE is not good when outliers are present.

**Mean Absolute Error (MAE):**

$$MAE = \sum_{n=1}^N |y_n - f(\mathbf{x}_n)|$$

**Huber loss**

$$Huber = \begin{cases} \frac{1}{2} z^2 & , |z| \leq \delta \\ \delta |z| - \frac{1}{2} \delta^2 & , |z| > \delta \end{cases}$$

- Huber loss is convex, differentiable, and also robust to outliers hard to set  $\delta$ .

**Tukey's bisquare loss**

$$L(z) = \begin{cases} z(\delta^2 - z^2)^2 & , |z| < \delta \\ 0 & , |z| \geq \delta \end{cases}$$

Tukey's loss is non-convex, non-differentiable, but robust to outliers.

**Hinge loss**

$$Hinge = [1 - y_n f(\mathbf{x}_n)]_+$$

**Logistic loss**

$$Logistic = \log(1 - \exp(y_n f(\mathbf{x}_n)))$$

## 3 Regression

- **Data** consists of  $N$  pairs  $(y_n, \mathbf{x}_n)$ 
  1.  $y_n$  the  $n$ 'th output
  2.  $\mathbf{x}_n$  is a vector of  $D$  inputs
- **Prediction:** predict the output for a new input vector.
- **Interpretation:** understand the effect of inputs on output.
- **Outliers** are data that are far away from most of the other examples.

### 3.1 Linear Regression

- Model that assume linear relationship between inputs and the output.
$$y_n \equiv f(\mathbf{x}_n) \\ =: \beta_0 + \beta_1 x_{n1} + \dots \\ = \mathbf{x}_n^T \boldsymbol{\beta}$$
with  $\boldsymbol{\beta}$  the parameters of the model.
- Variance grows only linearly with dimensionality

### 3.2 Gradient Descent

- Gradient descent uses only first-order information and takes steps in the direction of the gradient
- Given a cost function  $\mathcal{L}(\boldsymbol{\beta})$  we wish to find  $\boldsymbol{\beta}$  that minimizes the cost:

$$\min_{\boldsymbol{\beta}} \mathcal{L}(\boldsymbol{\beta})$$

#### 3.2.1 Grid search

- Compute the cost over a grid of  $M$  points to find the minimum
- Exponential Complexity  $O(NDM^D)$
- Hard to find a good range of values

### 3.3 Batch Gradient Descent

- Take steps in the opposite direction of the gradient

$$\boldsymbol{\beta}^{(k+1)} \leftarrow \boldsymbol{\beta}^{(k)} - \alpha \frac{d\mathcal{L}(\boldsymbol{\beta}^{(k)})}{d\boldsymbol{\beta}}$$

- with  $\alpha > 0$  the learning rate.
- With  $\alpha$  too big, method might diverge. With  $\alpha$  too small, convergence is slow.

### 3.4 Gradients for MSE

$$\tilde{\mathbf{X}} = [\mathbf{1} \quad \mathbf{X}]$$

- We define the error vector  $\mathbf{e}$ :

$$\mathbf{e} = \mathbf{y} - \tilde{\mathbf{X}} \boldsymbol{\beta}$$

- and MSE as follows:

$$\mathcal{L}(\boldsymbol{\beta}) = \frac{1}{2N} \sum_{n=1}^N (y_n - \tilde{\mathbf{x}}_n^T \boldsymbol{\beta})^2 = \frac{1}{2N} \mathbf{e}^T \mathbf{e}$$

- then the gradient is given by

$$\frac{d\mathcal{L}}{d\boldsymbol{\beta}} = -\frac{1}{N} \tilde{\mathbf{X}}^T \mathbf{e}$$

- Optimality conditions:

1. *necessary:* gradient equal to zero:

$$\frac{d\mathcal{L}(\boldsymbol{\beta}^*)}{d\boldsymbol{\beta}} = 0$$

2. *sufficient:* Hessian matrix is positive definite:

$$\mathbf{H}(\boldsymbol{\beta}^*) = \frac{d^2 \mathcal{L}(\boldsymbol{\beta}^*)}{d\boldsymbol{\beta} d\boldsymbol{\beta}^T}$$

- Very sensitive to illconditioning. Therefore, always normalize your feature otherwise step-size selection is difficult since different directions might move at different speed.
- *Complexity:*  $O(NDI)$  with  $I$  the number of iterations

### 3.5 Least Squares

- In some cases, we can compute the minimum of the cost function analytically.
- use the first optimality conditions:

$$\frac{d\mathcal{L}}{d\boldsymbol{\beta}} = 0 \Rightarrow \tilde{\mathbf{X}}^T \mathbf{e} = \tilde{\mathbf{X}}^T (\mathbf{y} - \tilde{\mathbf{X}} \boldsymbol{\beta}) = 0$$

- When  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  is invertible, we have the closed-form expression

$$\boldsymbol{\beta}^* = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

- thus we can predict values for a new  $\mathbf{x}_*$

$$\mathbf{y}_* = \tilde{\mathbf{x}}_*^T \boldsymbol{\beta}^* = \tilde{\mathbf{x}}_*^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

- The **Gram matrix**  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  is positive definite and is also invertible iff  $\tilde{\mathbf{X}}$  has full column rank.
- *Complexity:*  $O(ND^2 + D^3) \equiv O(ND^2)$
- $\tilde{\mathbf{X}}$  can be rank deficient when  $D > N$  or when the columns  $\tilde{\mathbf{x}}_i$  are nearly collinear. In this case, the matrix is ill-conditioned, leading to numerical issues.

### 3.6 Maximum Likelihood

- Let define our mistakes  $\epsilon_n \sim \mathcal{N}(0, \sigma^2)$ .

$$\rightarrow y_n = \tilde{\mathbf{x}}_n^T \boldsymbol{\beta} + \epsilon_n$$

- Another way of expressing this:

$$p(\mathbf{y}|\tilde{\mathbf{X}}, \boldsymbol{\beta}) = \prod_{n=1}^N p(y_n|\tilde{\mathbf{x}}_n, \boldsymbol{\beta}) \\ = \prod_{n=1}^N \mathcal{N}(y_n|\tilde{\mathbf{x}}_n^T \boldsymbol{\beta}, \sigma^2)$$

which defines the likelihood of observing  $\mathbf{y}$  given  $\tilde{\mathbf{X}}$  and  $\boldsymbol{\beta}$

- Define cost with log-likelihood

$$\mathcal{L}_{lik}(\boldsymbol{\beta}) = \log p(\mathbf{y}|\tilde{\mathbf{X}}, \boldsymbol{\beta})$$

$$= -\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \tilde{\mathbf{x}}_n^T \boldsymbol{\beta})^2 + cnst$$

- Maximum likelihood estimator (MLE) gives another way to design cost functions

$$\argmin_{\boldsymbol{\beta}} \mathcal{L}_{mse}(\boldsymbol{\beta}) = \argmin_{\boldsymbol{\beta}} \mathcal{L}_{lik}(\boldsymbol{\beta})$$

- MLE can also be interpreted as finding the model under which the observed data is most likely to have been generated from.
- With Laplace distribution

$$p(y_n|\tilde{\mathbf{x}}_n, \boldsymbol{\beta}) = \frac{1}{2b} e^{-\frac{1}{b} |y_n - \tilde{\mathbf{x}}_n^T \boldsymbol{\beta}|}$$

$$\sum_n \log p(y_n|\tilde{\mathbf{x}}_n, \boldsymbol{\beta}) = \sum_n |y_n - \tilde{\mathbf{x}}_n^T \boldsymbol{\beta}| + cnst$$

### 3.7 Ridge Regression

- Linear models usually underfit. One way is to use nonlinear basis functions instead.

$$y_n = \beta_0 + \sum_{j=1}^M \beta_j \phi_j(\mathbf{x}_n) = \tilde{\boldsymbol{\phi}}(\mathbf{x}_n)^T \boldsymbol{\beta}$$

- This model is linear in  $\boldsymbol{\beta}$  but nonlinear in  $\mathbf{x}$ . Note that the dimensionality is now  $M$ , not  $D$ .
- Polynomial basis

$$\boldsymbol{\phi}(x_n) = [1, x_n, x_n^2, \dots, x_n^M]$$

- The least square solution becomes

$$\boldsymbol{\beta}_{lse}^* = (\tilde{\boldsymbol{\Phi}}^T \tilde{\boldsymbol{\Phi}})^{-1} \tilde{\boldsymbol{\Phi}}^T \mathbf{y}$$

- Complex models overfit easily. Thus we can choose simpler models by adding a **regularization term** which penalizes complex models

$$\min_{\boldsymbol{\beta}} \left( \mathcal{L}(\boldsymbol{\beta}) + \frac{\lambda}{2N} \sum_{j=1}^M \beta_j^2 \right)$$

$$\boldsymbol{\beta}^* = \argmin_{\boldsymbol{\beta}} \left( \frac{1}{2} (\mathbf{y} - \mathbf{X} \boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}) + \frac{\lambda}{2} \boldsymbol{\beta}^T \boldsymbol{\beta} \right)$$

- Note that  $\beta_0$  is not penalized.
- By differentiating and setting to zero we get

$$\boldsymbol{\beta}_{ridge} = (\tilde{\boldsymbol{\Phi}}^T \tilde{\boldsymbol{\Phi}} + \boldsymbol{\Lambda})^{-1} \tilde{\boldsymbol{\Phi}}^T \mathbf{y}$$

$$\boldsymbol{\Lambda} = \begin{bmatrix} 0 & \mathbf{0} \\ \mathbf{0} & \lambda \mathbf{I}_m \end{bmatrix}$$

- Ridge regression improves the condition number of the Gram matrix since the eigenvalues of  $(\tilde{\boldsymbol{\Phi}}^T \tilde{\boldsymbol{\Phi}} + \lambda \mathbf{I}_m)$  are at least  $\lambda$
- **Maximum-a-posteriori (MAP) estimator:**
  - Maximizes the product of the likelihood and the prior.

$$\boldsymbol{\beta}_{map} = \argmax_{\boldsymbol{\beta}} (p(\mathbf{y}|\mathbf{X}, \boldsymbol{\Lambda}) p(\boldsymbol{\beta}|\boldsymbol{\Sigma}))$$

- Assume  $\beta_0 = 0$

$$\boldsymbol{\beta}_{ridge} = \argmax_{\boldsymbol{\beta}} \left( \log \left[ \prod_{n=1}^N \mathcal{N}(y_n|\tilde{\mathbf{x}}_n^T \boldsymbol{\beta}, \Lambda) \right] \times \mathcal{N}(\boldsymbol{\beta}|\mathbf{0}, \mathbf{I}) \right)$$

- **Lasso regularizer** forces some  $\beta_i$  to be strictly 0 and therefore forces sparsity in the model.

$$\min_{\boldsymbol{\beta}} \frac{1}{2N} \sum_{n=1}^N (y_n - \tilde{\boldsymbol{\phi}}(\mathbf{x}_n)^T \boldsymbol{\beta})^2, \quad \text{such that } \sum_{i=1}^M |\beta_i| \leq \tau$$

### 3.8 Cross-Validation

- We should choose  $\lambda$  to minimize the mistakes that will be made in the future.
- We split the data into train and validation sets and we pretend that the validation set is the future data. We fit our model on the training set and compute a prediction-error on the validation set. This gives us an *estimate* of the *generalization error*.
- **K-fold cross validation** randomly partition the data into  $K$  groups. We train on  $K-1$  groups and test on the remaining group. We repeat this until we have tested on all  $K$  sets. We then average the results.
- Cross-validation returns an unbiased estimate of the generalization error and its variance.

### 3.9 Bias-Variance decomposition

- The expected test error can be expressed as the sum of two terms
  - **Squared bias:** The average *shift* of the predictions
  - **Variance:** measure how data points vary around their average.

$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise}$$

- Both model bias and estimation bias are important
- Ridge regression increases estimation bias while reducing variance
- Increasing model complexity increases test error

Small  $\lambda \rightarrow$  low bias but large variance

Large  $\lambda \rightarrow$  large bias but low variance

### 3.10 Logistic Regression

- **Classification** relates input variables  $\mathbf{x}$  to discrete output variable  $y$
- **Binary classifier:** we use  $y = 0$  for  $\mathbf{C}_1$  and  $y = 1$  for  $\mathbf{C}_2$ .
- Can use least-squares to predict  $\hat{y}_*$

$$\hat{y} = \begin{cases} \mathbf{C}_1 & \hat{y}_* < 0.5 \\ \mathbf{C}_2 & \hat{y}_* \geq 0.5 \end{cases}$$

$$\sigma(x) = \frac{\exp(x)}{1 + \exp(x)}$$

$$p(y_n = \mathbf{C}_1|\mathbf{x}_n) = \sigma(\tilde{\mathbf{x}}_n^T \boldsymbol{\beta})$$

$$p(y_n = \mathbf{C}_2|\mathbf{x}_n) = 1 - \sigma(\tilde{\mathbf{x}}_n^T \boldsymbol{\beta})$$

- The probabilistic model:

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\beta}) = \prod_{n=1}^N \sigma(\tilde{\mathbf{x}}_n^T \boldsymbol{\beta})^{y_n} (1 - \sigma(\tilde{\mathbf{x}}_n^T \boldsymbol{\beta}))^{1-y_n}$$

- The log-likelihood:

$$\mathcal{L}_{mle}(\boldsymbol{\beta}) = \sum_{n=1}^N (y_n \tilde{\mathbf{x}}_n^T \boldsymbol{\beta} - \log(1 + \exp(\tilde{\mathbf{x}}_n^T \boldsymbol{\beta})))$$

- We can use the fact that

$$\frac{d}{dx} \log(1 + \exp(x)) = \sigma(x)$$

- Gradient of the log-likelihood

$$\mathbf{g} = \frac{d\mathcal{L}}{d\boldsymbol{\beta}} = \sum_{n=1}^N (\tilde{\mathbf{x}}_n y_n - \tilde{\mathbf{x}}_n \sigma(\tilde{\mathbf{x}}_n^T \boldsymbol{\beta}))$$

$$= -\tilde{\mathbf{X}}^T [\sigma(\tilde{\mathbf{X}} \boldsymbol{\beta}) - \mathbf{y}]$$

- The negative of the log-likelihood  $-\mathcal{L}_{mle}(\boldsymbol{\beta})$  is convex

- **Hessian** of the log-likelihood

$$\begin{aligned} & \text{We know that} \\ & \frac{d\sigma(t)}{dt} = \sigma(t)(1 - \sigma(t)) \end{aligned}$$

- Hessian is the derivative of the gradient

$$\begin{aligned} \mathbf{H}_{\boldsymbol{\beta}} &= -\frac{d\mathbf{g}(\boldsymbol{\beta})}{d\boldsymbol{\beta}^T} = \sum_{n=1}^N \frac{d}{d\boldsymbol{\beta}^T} \sigma(\tilde{\mathbf{x}}_n^T \boldsymbol{\beta}) \tilde{\mathbf{x}}_n \\ &= \sum_{n=1}^N \tilde{\mathbf{x}}_n \sigma(\tilde{\mathbf{x}}_n^T \boldsymbol{\beta}) (1 - \sigma(\tilde{\mathbf{x}}_n^T \boldsymbol{\beta})) \tilde{\mathbf{x}}_n^T \\ &= \tilde{\mathbf{X}}^T \mathbf{S} \tilde{\mathbf{X}} \end{aligned}$$

where  $\mathbf{S}$  is a  $N \times N$  diagonal matrix with diagonals

$$S_{nn} = \sigma(\tilde{\mathbf{x}}_n^T \boldsymbol{\beta}) (1 - \sigma(\tilde{\mathbf{x}}_n^T \boldsymbol{\beta}))$$

- The negative of the log-likelihood is not strictly convex.

### - Newton's Method

- Uses second-order information and takes steps in the direction that minimizes a quadratic approximation

$$\mathcal{L}(\boldsymbol{\beta}) = \mathcal{L}(\boldsymbol{\beta}^{(k)}) + \mathbf{g}_k^T (\boldsymbol{\beta} - \boldsymbol{\beta}^{(k)}) + (\boldsymbol{\beta} - \boldsymbol{\beta}^{(k)})^T \mathbf{H}_k (\boldsymbol{\beta} - \boldsymbol{\beta}^{(k)})$$

and it's minimum is at

$$\boldsymbol{\beta}^{k+1} = \boldsymbol{\beta}^{(k)} - \alpha_k \mathbf{H}_k^{-1} \mathbf{g}_k$$

where  $\mathbf{g}_k$  is the gradient and  $\alpha_k$  the learning rate.

- Complexity:  $O((ND^2 + D^3)I)$

### - Penalized Logistic Regression

$$\min_{\boldsymbol{\beta}} \left( -\sum_{n=1}^N \log p(y_n|\tilde{\mathbf{x}}_n^T \boldsymbol{\beta}) + \lambda \sum_{d=1}^D \beta_d^2 \right)$$

## 4 Generalized Linear Model

### - Exponential family distribution

$$p(\mathbf{y}|\boldsymbol{\eta}) = \frac{h(\mathbf{y})}{Z} \exp(\boldsymbol{\eta}^T \boldsymbol{\phi}(\mathbf{y}) - A(\boldsymbol{\eta}))$$

- Bernoulli distribution

$$p(y|\mu) = \mu^y (1 - \mu)^{1-y}$$

$$= \exp(y \log(\frac{\mu}{1-\mu}) + \log(1-\mu))$$

- there is a relationship between  $\boldsymbol{\eta}$  and  $\mu$  through the **link function**

$$\boldsymbol{\eta} = \log\left(\frac{\mu}{1-\mu}\right) \Leftrightarrow \mu = \frac{e^{\boldsymbol{\eta}}}{1 + e^{\boldsymbol{\eta}}}$$

- Note that  $\mu$  is the mean parameter of  $y$
- Relationship between the mean  $\mu$  and  $\boldsymbol{\eta}$  is defined using a link function  $g$

$$\boldsymbol{\eta} = \mathbf{g}(\boldsymbol{\mu}) \Leftrightarrow \boldsymbol{\mu} = \mathbf{g}^{-1}(\boldsymbol{\eta})$$

- First and second derivatives of  $A(\boldsymbol{\eta})$  are related to the mean and the variance

$$\frac{dA(\boldsymbol{\eta})}{d\boldsymbol{\eta}} = E[\boldsymbol{\phi}(\boldsymbol{\eta})], \quad \frac{d^2 A(\boldsymbol{\eta})}{d$$

- We obtain the solution

$$\frac{d\mathcal{L}}{d\beta} = \tilde{\mathbf{X}}^T [\mathbf{g}^{-1}(\eta) - \phi(\mathbf{y})]$$

## 5 k-Nearest Neighbor (k-NN)

- The k-NN prediction for  $\mathbf{x}$  is

$$f(\mathbf{x}) = \frac{1}{k} \sum_{\mathbf{x}_n \in nbh_k(\mathbf{x})} y_n$$

- where  $nbh_k(\mathbf{x})$  is the neighborhood of  $\mathbf{x}$  defined by the  $k$  closest points  $\mathbf{x}_n$  in the training data
- **Curse of dimensionality:** Generalizing correctly becomes exponentially harder as the dimensionality grows.
- Gathering more inputs variables may be a bad thing

## 6 Kernel Ridge Regression

- The following is true for ridge regression

$$\beta = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^T \mathbf{y} \quad (1)$$

$$= \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I}_N)^{-1} \mathbf{y} \quad (2)$$

$$= \mathbf{X}^T \alpha$$

- Complexity of computing  $\beta$

1.  $O(D^2 N + D^3)$
2.  $O(N^2 D + N^3)$

- Thus we have

$$\beta = \sum_{n=1}^D \alpha_n \mathbf{x}_n, \quad \mathbf{y} = \sum_{d=1}^D \beta_d \tilde{\mathbf{x}}_d$$

- with  $\mathbf{x}_n$  the rows of  $\mathbf{X}$  and  $\tilde{\mathbf{x}}_d$  the columns of  $\mathbf{X}$
- The representer theorem allows us to write an equivalent optimization problem in terms of  $\alpha$ .

$$\alpha = \arg \max_{\alpha} \left( -\frac{1}{2} \alpha (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I}_N)^T \alpha + \alpha^T \mathbf{y} \right)$$

- $\mathbf{X} \mathbf{X}^T$  is called the **kernel matrix**
- Kernelized ridge regression might be computationally more efficient in some cases.
- **Kernel trick:**
  - We can work directly with  $\mathbf{K} = \mathbf{X} \mathbf{X}^T$  and never have to worry about  $\mathbf{X}$
  - Using the basis function  $\phi(\mathbf{x})$ , we do not need to specify it explicitly, since we can work directly with  $\mathbf{K} = \phi(\mathbf{x})$
  - We will use a kernel function  $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$
  - The evaluation of a kernel is usually faster with  $k$  than with  $\phi$
- **Radial Basis function kernel (RBF)**

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^T(\mathbf{x} - \mathbf{x}')\right)$$

- Properties of a kernel to ensure the existence of a corresponding  $\phi$ :
  - $\mathbf{K}$  should be symmetric:  $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$
  - $\mathbf{K}$  should be positive semidefinite.

## 7 Support Vector Machine

- **Hinge loss**

$$[t]_+ = \max(0, t)$$

- Solution to the dual problem is sparse and non-zero entries will be our **support vectors**.
- Assume  $y_n \in \{-1, 1\}$
- SVM optimizes the following cost

$$\min_{\beta} \sum_{n=1}^N [1 - y_n \tilde{\phi}_n^T \beta]_+ + \frac{\lambda}{2} \sum_{j=1}^M \beta_j^2$$

- The minimum doesn't change with a rescaling of  $\beta$
- **Duality:**

- Hard to minimize  $g(\beta)$  so we define  $g(\beta) = \max_{\alpha} G(\beta, \alpha)$

- we use the property that  $C[v_n]_+ = \max(0, C v_n) = \max_{\alpha_n \in [0, C]} \alpha_n v_n$

- We can rewrite the problem as

$$\min_{\beta} \max_{\alpha \in [0, C]^N} \sum_{n=1}^N \alpha_n (1 - y_n \tilde{\phi}_n^T \beta) + \frac{1}{2} \sum_{j=1}^M \beta_j^2$$

- This is differentiable, convex in  $\beta$  and concave in  $\alpha$

- **Minimax theorem:**

$$\min_{\beta} \max_{\alpha} G(\beta, \alpha) = \max_{\beta} \min_{\alpha} G(\beta, \alpha)$$

- because  $G$  is convex in  $\beta$  and concave in  $\alpha$ .
- Derivative w.r.t.  $\beta$ :

$$\frac{dG}{d\beta} = - \left( \sum_{n=1}^N \alpha y_n \tilde{\phi}_n \right) + \begin{bmatrix} 0 \\ \beta_{1:M} \end{bmatrix}$$

- Equating this to 0, we get:

$$\beta_{1:M}^* = \sum_{n=1}^N \alpha_n y_n \tilde{\phi}_n = \Phi^T \text{diag}(\mathbf{y}) \alpha$$

$$\alpha^T \mathbf{y} = 0$$

- Plugging  $\beta^*$  back in the dual problem

$$\max_{\alpha \in [0, C]^N} \alpha^T \mathbf{1} - \frac{1}{2} \alpha^T \mathbf{Y} \Phi \Phi^T \mathbf{Y} \alpha$$

- This is a differentiable least-squares problem. Optimization is easy using Sequential Minimal Optimization. It is also naturally kernelized with  $\mathbf{K} = \Phi \Phi^T$
- The solution  $\alpha$  is sparse and is non-zero only for the training examples that are instrumental in determining the decision boundary.

## 8 K-means

- **Unsupervised learning:** Represent particular input patterns in a way that reflects the statistical structure of the overall collections of input patterns.
- **Cluster** are groups of points whose inter-point distances are small compared to the distances outside the cluster.

$$\min_{\mathbf{r}, \mu} \sum_{k=1}^K \sum_{n=1}^N r_{nk} \|\mathbf{x}_n - \mu_k\|_2^2$$

- such that  $r_{nk} \in \{0, 1\}$  and  $\sum_{k=1}^K r_{nk} = 1$
- K-means algorithm: Initialize  $\mu_k$ , then iterate

1. For all  $n$ , compute  $r_n$  given  $\mu$

$$r_{nk} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\mathbf{x}_n - \mu_j\|_2^2 \\ 0 & \text{otherwise} \end{cases}$$

2. For all  $k$ , compute  $\mu_k$  given  $\mathbf{r}$

$$\mu_k = \frac{\sum_{n=1}^N r_{nk} \mathbf{x}_n}{\sum_{n=1}^N r_{nk}}$$

- A good initialization procedure is to choose the prototypes to be equal to a random subset of  $K$  data points.
- Probabilistic model

$$p(\mathbf{r}, \mu) = \prod_{n=1}^N \prod_{k=1}^K [\mathcal{N}(\mathbf{x}_n | \mu_k, \mathbf{I})]^{r_{nk}}$$

- Computation can be heavy, each example can belong to only on cluster and clusters have to be spherical.

## 9 Gaussian Mixture Models

- Clusters can be spherical using a full covariance matrix instead of isotropic covariance.

$$p(\mathbf{X} | \mu, \Sigma, \mathbf{r}) = \prod_{n=1}^N \prod_{k=1}^K [\mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)]^{r_{nk}}$$

- **Soft-clustering:** Points can belong to several cluster by defining  $r_n$  to be a random variable.

$$p(r_{nk} = 1) = \pi_k \text{ where } \pi_k > 0, \forall k$$

$$\sum_{k=1}^K \pi_k = 1$$

- Joint distribution of Gaussian mixture model

$$p(\mathbf{X}, \mathbf{r} | \mu, \Sigma, \pi) = \prod_{n=1}^N [p(\mathbf{x}_n | r_n, \mu, \Sigma) p(r_n | \pi)]$$

$$= \left[ \prod_{k=1}^K [(\mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k))^{r_{nk}}] \prod_{k=1}^K [\pi^{r_{nk}}] \right]$$

$$\text{Joint} = \text{Likelihood} \times \text{Prior}$$

- $r_n$  are called *latent* unobserved variables
- Unknown parameters are given by  $\theta = \{\mu, \Sigma, \pi\}$

- We get the **marginal likelihood** by marginalizing  $r_n$  out from the likelihood

$$p(\mathbf{x}_n | \theta) = \sum_{k=1}^K p(\mathbf{x}_n, r_n = k | \theta)$$

$$= \sum_{k=1}^K p(r_n = k | \theta) p(\mathbf{x}_n | r_n = k, \theta)$$

$$= \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)$$

- Without a latent variable model, number of parameters grow at rate  $O(N)$
- After marginalization, the growth is reduced to  $O(D^2 K)$
- To get maximum likelihood estimate of  $\theta$ , we maximize

$$\max_{\theta} \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \sigma_k)$$

## 10 Expectation Maximization Algorithm

- [ALGORITHM] Start with  $\theta^{(1)}$  and iterate
  1. *Expectation step:* Compute a lower bound to the cost such that it is tight at the previous  $\theta^{(i)}$

$$\log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \geq \sum_{k=1}^K p_{kn} \log \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}$$

with equality when,

$$p_{kn} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n, \Sigma_k)}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}$$

2. *Maximization step:* Update  $\theta$   
 $\theta^{(i+1)} = \arg \max_{\theta} \mathcal{L}(\theta, \theta^{(i)})$

$$\mu_k^{i+1} = \frac{\sum_{n=1}^N p_{kn}^{(i)} \mathbf{x}_n}{\sum_{n=1}^N p_{kn}^{(i)}}$$

$$\Sigma_k^{(i+1)} = \frac{\sum_{n=1}^N p_{kn}^{(i)} (\mathbf{x}_n - \mu_k^{(i+1)}) (\mathbf{x}_n - \mu_k^{(i+1)})^T}{\sum_{n=1}^N p_{kn}^{(i)}}$$

$$\pi_k^{(i+1)} = \frac{1}{N} \sum_{n=1}^N p_{kn}^{(i)}$$

- If the covariance is diagonal, then we have K-means.

## 11 Matrix factorization

- We have  $D$  movies and  $N$  users
- $\mathbf{X}$  is a matrix  $D \times N$  with  $x_{dn}$  the rating of  $n$ 'th user for  $d$ 'th movie.
- We project data vectors  $\mathbf{x}_n$  to a smaller dimension  $\mathbf{z}_n \in \mathbb{R}^M$
- We have now 2 latent variables:
  - $\mathbf{Z}$  a  $N \times M$  matrix that gives features for the users
  - $\mathbf{W}$  a  $D \times M$  matrix that gives features for the movies

$$x_{dn} \approx \mathbf{w}_d^T \mathbf{z}_n$$

- We can add a regularizer and minimize the following cost:

$$\mathcal{L}(\mathbf{W}, \mathbf{Z}) = \frac{1}{2} \sum_{n=1}^N \sum_{d=1}^D (x_{dn} - \mathbf{w}_d^T \mathbf{z}_n)^2 + \frac{\lambda_w}{2} \sum_{d=1}^D \|\mathbf{w}_d\|^2 + \frac{\lambda_z}{2} \sum_{n=1}^N \|\mathbf{z}_n\|^2$$

- We can use coordinate descent algorithm, by first minimizing w.r.t.  $\mathbf{Z}$  given  $\mathbf{W}$  and then minimizing  $\mathbf{W}$  given  $\mathbf{Z}$ . This is called **Alternating least-squares (ALS)**:

$$\mathbf{Z}^T \leftarrow (\mathbf{W}^T \mathbf{W} + \lambda_z \mathbf{I}_M)^{-1} \mathbf{W}^T \mathbf{X}$$

$$\mathbf{W}^T \leftarrow (\mathbf{Z}^T \mathbf{Z} + \lambda_w \mathbf{I}_M)^{-1} \mathbf{Z}^T \mathbf{X}^T$$

- Complexity:  $O(DNM^2 + NM^3) \rightarrow O(DNM^2)$
- Probabilistic model

$$\prod_{n=1}^N \prod_{d \in \mathcal{O}_n} \mathcal{N}(x_{dn} | \mathbf{w}_d^T \mathbf{z}_n, I) \times \prod_{n=1}^N \mathcal{N}(\mathbf{z}_n | 0, \frac{1}{\lambda_z} I)$$

- Since many ratings are missing we cannot

- normalize the data. A solution is to add offset terms:

$$-\frac{1}{2} \sum_{n=1}^N \sum_{d \in \mathcal{O}_n} (x_{dn} - \mathbf{w}_d^T \mathbf{z}_n - w_{0d} - z_{0n} - \mu)^2$$

## 12 Singular Value Decomposition

- Matrix factorization method

$$\mathbf{X} = \mathbf{U} \mathbf{S} \mathbf{V}^T$$

- $\mathbf{U}$  is an  $D \times D$  matrix
- $\mathbf{V}$  is an  $N \times N$  matrix
- $\mathbf{S}$  is a non-negative diagonal matrix of size  $D \times N$  which are called **singular values** appearing in a descending order.
- Columns of  $\mathbf{U}$  and  $\mathbf{V}$  are the left and right **singular vectors** respectively.
- Assuming  $D < N$  we have

$$\mathbf{X} = \sum_{d=1}^D s_d \mathbf{u}_d \mathbf{v}_d^T$$

- This tells you about the spectrum of  $\mathbf{X}$  where higher singular vectors contain the *low-frequency information* and lower singular values contain the *high-frequency information*.
- Let's now truncate these matrices

$$\mathbf{X} \approx \mathbf{U}_{tr} \mathbf{S}_{tr} \mathbf{V}_{tr}^*$$

$$\mathbf{T}_{tr} \approx \mathbf{U}_{tr} \mathbf{S}_{tr}$$

$$\mathbf{T}_{te} \approx \mathbf{X}_{te} \mathbf{V}_{tr}$$

- $\mathbf{p}_{kn}$  with  $\mathbf{T}_{tr}$  the reduced feature set of  $\mathbf{X}$  and  $\mathbf{T}_{te}$  the reduced feature set of  $\mathbf{X}_{te}$

## 13 Principal Component Analysis

- PCA is a dimensionality reduction method and a method to decorrelate the data

$$\mathbf{x} \approx \tilde{\mathbf{x}} = \mathbf{W} \mathbf{z}^T$$

- such that columns of  $\mathbf{W}$  are orthogonal.
- If the data is zero mean

$$\mathbf{C} = \frac{1}{N} \mathbf{X} \mathbf{X}^T \approx \mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{S}^2 \mathbf{U}^T$$

$$\Rightarrow \mathbf{U}^T \mathbf{X} \mathbf{X}^T \mathbf{U} = \mathbf{U}^T \mathbf{U} \mathbf{S}^2 \mathbf{U}^T \mathbf{U} = \mathbf{S}^2$$

- The columns of matrix  $\mathbf{U}$  are called the **principal components** and they decorrelate the covariance matrix.
- Using SVD, we can compute the matrices in the following way

$$\mathbf{W} = \mathbf{U} \mathbf{S}^{1/2}$$

$$\mathbf{Z} = \mathbf{V} \mathbf{S}^{1/2}$$

## 14 Belief Propagation

- the goal is to learn *inference of the latent variables using belief propagation*
- Given a directed acyclic graph  $G$  and parameters  $\theta$ , a **Bayesian network** defines the joint distribution as follows

$$p_{\theta}(\mathbf{x}) = \prod_{k=1}^K p_{\theta}(x_k | \text{parents}_k)$$

- We can reduce the complexity by using the structure of the problem and the bayesian rule.

## 15 Multi-Layer Perceptron (MLP)

- Known as **feed-forward neural network**

- $\mathbf{z}_n^{(k)}$  is the  $k$ 'th hidden vector
- $\mathbf{a}_n^{(k)} = (\beta_n^{(k)})^T \mathbf{z}_n^{(k-1)}$ ,  $\mathbf{z}_n^{(k)} = h(\mathbf{a}_n^{(k)})$

- For the first layer we have  $\mathbf{z}_n^{(0)} = \mathbf{x}_n$
- For the last layer, we use a link function to map  $\mathbf{z}_n^{(K-1)}$  to the output  $\mathbf{y}_n$

- A 1-layer MLP is simply a generalization of linear/logistic regression
- $\mathbf{B}^{(k)}$  a matrix with rows  $(\beta_m^{(k)})^T$
- $\mathbf{a}_n^{(k)} = \mathbf{B}^{(k)} \mathbf{z}_n^{(k-1)}$ ,  $\mathbf{z}_n^{(k)} = h(\mathbf{a}_n^{(k)})$

- Thus we have the input-output relationship  $\mathbf{y}_n = g((\beta_n^{(K-1)})^T * h(\mathbf{B}^{(K-2)} * h(\dots * h(\mathbf{B}^{(1)} * \mathbf{x}_n)))$  with  $g()$  the link function

- We learn parameters  $\mathbf{B}$  using stochastic gradient-descent
- Frequently used transfer function

$$g(a) = \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$

- **Backpropagation** is a technic to compute the gradient in time linear in the number of training points and the number of weights.

## 16 Gaussian Process (GP)

- GP are a family of statistical distributions in which time plays a role and for which any finite linear combination of samples has a joint Gaussian distribution.
- They can be completely defined by their second-order statistics which means that if they have mean zero, then defining the covariance function completely defines the process behaviour.
- Let us place a probabilistic prior shape on the approximation of a function.
- A GP process defines a prior over function  $f$

$$p(f|X) = \mathcal{N}(f|0, K(X))$$

- $K(X)$  defines shape and prior knowledge about our problem

$$p(f|\mathbf{y}, \mathbf{X}_*, \mathbf{X}) \sim \mathcal{N}(\mu', \sigma'^2)$$

$$\mu' = K(\mathbf{X}_*, \mathbf{X}) [K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{y}$$

$$\sigma' = K(\mathbf{X}_*, \mathbf{X}_*) - K(\mathbf{X}_*, \mathbf{X}) [K(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}] K(\mathbf{X}, \mathbf{X}_*)$$

- with  $\sigma_n$  the variance of the noise
- RBF kernel

$$k(\mathbf{x}_n, \mathbf{x}_m) = e^{-\|\mathbf{x}_n - \mathbf{x}_m\|^2 / L^2}$$

- Quadratic kernel

$$k(\mathbf{x}_n, \mathbf{x}_m) = (1 + \mathbf{x}_n^T \mathbf{x}_m)^2$$

## 17 Decision Trees (DT)

- A decision tree is a structure in which each internal node represents a test on an attribute and each branch represents the outcome of the test and each leaf represents a class label.
- Fast to train and fast to make predictions
- Efficient for very high dimensional feature spaces and very large amounts of training data
- Lack of smoothness and high variance (overfitting)
- Goal: find a split  $(k, \tau)$  that minimizes an impurity measure at the leaves
  - Find best feature to split on
  - Find best threshold

## 18 Random Forests (RF)

- RF correct the overfitting bad "habit" of DTs.
- Training: Learn  $M$  trees on different subsets (random) of training data
- Prediction: Average of prediction of each tree
- Variance of the model averaging:

$$z_1 = f_1 \Rightarrow z_M = \frac{1}{M} \sum_{i=1}^M f_i$$

$$V(z_1) = \sigma^2 \Rightarrow V(z_M) = \frac{1}{M} \sigma^2 + \rho \frac{M-1}{M} \sigma^2$$

- Variance reduction ratio:

$$\frac{V(z_1)}{V(z_M)} = \frac{M}{1 + \rho(M-1)}$$

- 2 techniques for decorrelating trees:

- **Bagging:** Randomize training data
- Randomized feature selection