

Linear Regression: Find Optimal Model

• You want to minimize $E_{in}(\underline{w}) = (\underline{\underline{X}}\underline{w} - \underline{y})^T (\underline{\underline{X}}\underline{w} - \underline{y})$ with respect to \underline{w}

$$\nabla E_{in}(\underline{\boldsymbol{w}}^*) = \underline{\boldsymbol{0}}$$

$$\frac{2}{N}\underline{\underline{\boldsymbol{X}}}^T(\underline{\underline{\boldsymbol{X}}}\underline{\boldsymbol{w}}^* - \underline{\boldsymbol{y}}) = \underline{\boldsymbol{0}}$$

$$\underline{\underline{\boldsymbol{X}}}^T\underline{\underline{\boldsymbol{X}}}\underline{\boldsymbol{w}}^* = \underline{\underline{\boldsymbol{X}}}^T\underline{\boldsymbol{y}}$$

• If the square matrix $\underline{\mathbf{X}}^T\underline{\mathbf{X}}$ is invertible:

$$\underline{\boldsymbol{w}}^* = (\underline{\underline{\boldsymbol{X}}}^T \underline{\underline{\boldsymbol{X}}})^{-1} \underline{\underline{\boldsymbol{X}}}^T \underline{\boldsymbol{y}} = \underline{\underline{\boldsymbol{X}}}^\dagger \underline{\boldsymbol{y}}$$

- The matrix $\underline{\underline{X}}^\dagger \triangleq (\underline{\underline{X}}^T\underline{\underline{X}})^{-1}\underline{\underline{X}}^T$ is called **pseudo-inverse**
 - It generalizes the inverse for non-square matrices, in fact:
 - $\underline{X}^{\dagger}\underline{X} = \underline{I}$
 - If \underline{X} is square and invertible, then $\underline{X}^{\dagger} = \underline{X}^{-1}$



Complexity of One-Step Learning

- Learning with the pseudo-inverse is one-step learning
 - Contrasts with iterative methods, e.g., gradient descent
- Inverting a square matrix of size P is related to the number of parameters, not examples N
 - Complexity of one-step learning is $O(P^3)$



Linear Models Are Linear in What?

- A model is linear when the signal $s = \sum_{i=0}^{P} w_i x_i = \underline{\mathbf{w}}^T \underline{\mathbf{x}}$ is linear with variables
 - Unknown variables: weights w_i
 - Inputs x_i are fixed
- Applying a **non-linear transform to inputs** $z_i = \Phi(x_i)$ keeps the model linear, e.g.,
 - Positive/negative part (e.g., $z_i = RELU(x_i), RELU(-x_i)$)
 - Waterfall (conditioning model to different feature ranges)
 - Thresholding (e.g., $z_i = \min(x_i, T)$)
 - Indicator variables (e.g., $z_i = I(x_i > 0)$)
 - Winsorizing (replace outliers with a large constant value)
- Applying a non-linear transform to weights z_i = Φ(w_i) makes the model non-linear



Non-Linear Transformations in Linear Models

Transform variables

- Use $\Phi: \mathcal{X} \to \mathcal{Z}$
- Transform each point $\underline{\mathbf{x}}_n \in \mathcal{X} = \mathbb{R}^d$ into a point in feature space $\mathbf{z}_n = \Phi(\mathbf{x}_n) \in \mathcal{Z} = \mathbb{R}^{\tilde{d}}$ with $d \neq \tilde{d}$
- Learn
 - Learn linear model in \mathcal{Z} , obtaining $\underline{\tilde{w}}$ for separating hyperplane
- Predict
 - Evaluate model on new point in \mathcal{Z} :

$$y = \operatorname{sign}(\underline{\tilde{\boldsymbol{w}}}^T \Phi(\underline{\boldsymbol{x}})) \text{ or } y = \underline{\tilde{\boldsymbol{w}}}^T \Phi(\underline{\boldsymbol{x}})$$

- Compute decision boundary
 - In \mathcal{X} if Φ is invertible; or
 - ullet By classifying any $\underline{{\it x}} \in \mathcal{X}$ in \mathcal{Z}



- Perceptron
- Logistic Regression
- LDA, QDA
- Kernel Methods



Example of Classification Problems

- Binary classification problem
 - $y \in \{0, 1\}$
 - Typically assign 1 to what you want to detect
 - Email: spam, not_spam
 - Online transaction: fraudulent, valid
 - Tumor: malignant, benign
- Multi-class classification problem
 - $y \in \{0, 1, 2, ..., K\}$
 - Email tagging: work, family, friends
 - Medical diagnosis: not_ill, cold, flu
 - Weather: sunny, rainy, cloudy, snow



Linear Regression for Classification

- Use linear regression for classification
 - Transform outputs into $\{+1,-1\} \in \mathbb{R}$
 - Learn $\underline{\boldsymbol{w}}^T\underline{\boldsymbol{x}}_n \approx y_n = \pm 1$
 - Use sign $(\underline{w}^T \underline{x}_n)$ as model (perceptron)
- Not optimal: outliers influence fit due to square distance metric
 - Use weights from linear regression to initialize a learning algorithm for classification (e.g., PLA)



Perceptron Learning Algorithm (PLA)

- First machine learning algorithm discovered
- Algorithm
 - Training set $\mathcal{D} = \{(\underline{x}_1, y_1), ..., (\underline{x}_n, y_n)\}$
 - Initialize weights w
 - Random values
 - Use linear regression for classification as seed
 - Pick a misclassified point sign $(\underline{\boldsymbol{w}}^T\underline{\boldsymbol{x}}_i) \neq y_i$ from training set \mathcal{D}
 - Update weights: $\underline{\mathbf{w}}(t+1) = \underline{\mathbf{w}}(t) + y_i \underline{\mathbf{x}}_i$
 - Like stochastic gradient descent
 - Iterate until no misclassified points
- Algorithm converges (slowly) for linearly separable data
- Pocket version of PLA
 - Idea
 - Continuously update solution
 - Keep best solution "in the back pocket"
 - Have a solution if stopping after max iterations



Non-Linear Transformations for Classifications

- Classification problems have varying degrees of non-linear boundaries:
 - 1. Non-linearly separable data
 - E.g., + in center, in corners in a 2-feature scatter plot
 - 2. Mostly linear classes with few outliers
 - 3. Higher-order decision boundary
 - E.g., quadratic
 - 4. Non-linear data-feature relationship
 - E.g., variable threshold



- Perceptron
- Logistic Regression
- LDA, QDA
- Kernel Methods



Logistic Regression Is a Probabilistic Classifier

- Logistic regression learns:
 - The probability of each class Pr(y|x) given input x
 - Instead of predicting class y directly
- Parametric approach: assume $Pr(y|\underline{x};\underline{w})$ has a known functional form

$$\Pr(y = 1 | \underline{x}; \underline{w}) = \operatorname{logit}(\underline{w}^T \underline{x})$$

• Optimize parameters <u>w</u> using maximum likelihood

$$\underline{\mathbf{w}}^* = \operatorname{argmax}_{\mathbf{w}} \operatorname{Pr}(\mathcal{D}; \underline{\mathbf{w}})$$

• Predict by outputting class with highest probability

$$h_{\underline{\boldsymbol{w}}}(\underline{\boldsymbol{x}}) = \begin{cases} +1 & \Pr(y = 1 | \underline{\boldsymbol{x}}; \underline{\boldsymbol{w}}) \ge 0.5 \\ -1 & \Pr(y = 1 | \underline{\boldsymbol{x}}; \underline{\boldsymbol{w}}) < 0.5 \end{cases}$$



Logistic Regression: Example

- Assume y is:
 - The event y_i "patient with characteristics \mathbf{x} had heart attack"
 - Function of parameters \underline{x} (e.g., age, gender, diet)
- In data set \mathcal{D} :
 - No samples of Pr(y|x)
 - Have realizations:
 - "Patient with x₁ had a heart attack"
 - "Patient with x2 didn't"
 - •
- Learn $Pr(y|\underline{x})$
 - \bullet Find best parameters $\underline{\textbf{\textit{w}}}$ for logistic regression model to explain data \mathcal{D}



Logistic Function

- Aka "sigmoid"
- Logistic function logit(s) is defined as

$$\mathsf{logit}(s) \triangleq \frac{e^s}{1 + e^s} = \frac{1}{1 + e^{-s}}$$

- Varies in [0, 1]
- Crosses the origin at 0.5
- Asymptotes at 0 and 1
- It is a soft version of sign()



Logistic Regression vs Linear Classifier

- Functional form is similar
 - Logistic regression: $h(\mathbf{w}) = \text{logit}(\mathbf{w}^T \mathbf{x})$
 - Linear classifier (perceptron): $h(\underline{w}) = \text{sign}(\underline{w}^T \underline{x})$
- Difference in probabilistic interpretation and fitting method
 - Logistic regression** lacks samples of probability function to interpolate
 - Uses realizations of random variable
 - · Seeks model parameters maximizing data likelihood
 - Linear classification assumes class value is linear function of inputs



Error for Probabilistic Binary Classifiers

 For probabilistic binary classification, use log-probability error as point-wise error

$$e(h(\underline{x}), y) \triangleq -\log(\Pr(y = h(\underline{x})|\underline{x}; \underline{w}))$$

- Negate for positive errors: $log([0,1]) \in [-\infty,0)$
- Log probability generalizes 0-1 error
 - Case *y* = 1
 - Output $h(\underline{x})$ close to $1 \implies \text{probability } 1 \implies \log(1) = 0 \implies e(\cdot) = 0$
 - Output close to $0 \implies e() = -\log(0) \rightarrow +\infty$
 - Similar behavior for y = 0

One-Liner Error for Probabilistic Binary Classifiers

 Point-wise error for example (x, y) for probabilistic binary classifiers is defined as:

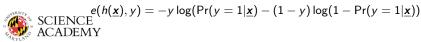
$$\begin{split} e(h(\underline{x}), y) &\triangleq -\log(\Pr(h(\underline{x}) = y | \underline{x})) \\ &= \begin{cases} -\log(\Pr(y = 1 | \underline{x})) & y = 1 \\ -\log(\Pr(y = 0 | \underline{x})) & y = 0 \end{cases} \\ &= \begin{cases} -\log(\Pr(y = 1 | \underline{x})) & y = 1 \\ -\log(1 - \Pr(y = 1 | \underline{x})) & y = 0 \end{cases} \end{split}$$

Any function of a binary variable:

$$y = \begin{cases} a & x = 1 \\ b & x = 0 \end{cases}$$

can be written as one-liner: $y = x \cdot a + (1 - x) \cdot b$

• Point-wise error can be written independently of $Pr(y = 1 | \underline{x})$:



One-Liner Error for Logistic Regression

The point-wise error for a binary classifier is:

$$e(h(\underline{x}), y) = -y \log(\Pr(y = 1|\underline{x}) - (1 - y) \log(1 - \Pr(y = 1|\underline{x}))$$

• Simplify further with logit function:

$$\begin{split} e(h(\underline{x}),y) &\triangleq -\log \Pr(h(\underline{x}) = y) \\ &\text{... a bit of math manipulation ...} \\ &= -\log \log \operatorname{it}(y\underline{\boldsymbol{w}}^T\underline{\boldsymbol{x}}) \qquad \text{since logit}(s) = \frac{1}{1+e^{-s}} \\ &= \log(1+\exp(-y\underline{\boldsymbol{w}}^T\underline{\boldsymbol{x}})) \end{split}$$

Point-wise error for logistic regression equals cross-entropy error



Cross-Entropy Error

Point-wise error for logistic regression has expression:

$$e(h(\underline{x}), y) = \log(1 + \exp(-y \cdot \underline{w}^T \underline{x}))$$

- It is called cross-entropy error
- Note: no before $log(\cdot)$ but before $y \cdot \underline{w}^T \underline{x}$
- Cross-entropy error generalizes 0-1 error
 - If $\underline{\mathbf{w}}^T \underline{\mathbf{x}}$ agrees with y in sign and $|\underline{\mathbf{w}}^T \underline{\mathbf{x}}|$ is large \implies error goes to 0
 - ullet If they disagree in sign \Longrightarrow error goes towards ∞
- Define in-sample error on training set as average of point-wise errors:

$$E_{in} \triangleq \frac{1}{N} \sum_{n} e(h(\underline{x}_n), y_n)$$

Fitting Logistic Regression

A plausible error measure of a hypothesis is based on likelihood of data

$$Pr(\mathcal{D}|h=f)$$

- "How likely is the data \mathcal{D} under the assumption that h = f?"
- "How likely is that the data $\mathcal D$ was generated by h?"
- Maximize likelihood $\mathcal D$ generated from logistic regression

$$\Pr(y=1|\underline{x};\underline{w})$$

 It can be proved that this is equivalent to minimizing in-sample error on training set using cross-entropy error



Fitting Logistic Regression (Optional)

• Find \underline{w} that maximizes likelihood for data set $\mathcal{D} = \{(\underline{x}_1, y_1), ..., (\underline{x}_N, y_N)\}$ generated by model $h(\underline{x})$:

$$\Pr(D|\underline{\boldsymbol{w}}) = \Pr(y_1 = h(\underline{\boldsymbol{x}}_1) \wedge ... \wedge y_N = h(\underline{\boldsymbol{x}}_N)) = \Pr(y_1 = y_1' \wedge ... \wedge y_N = y_N')$$

• Model form:

$$y' = h(\underline{x}) = \begin{cases} +1 & \text{if logit}(\underline{w}^T\underline{x}) > 0.5 \\ -1 & \text{otherwise} \end{cases}$$

Assuming independence among training examples:

$$Pr(D|\underline{\boldsymbol{w}}) = \prod_{i=1}^{N} Pr(y_i = y_i'|\underline{\boldsymbol{x}}_i)$$

• Fold y_n in expression:

• When
$$y_n = 1$$
, $Pr(y_n = y'_n) = logit(\underline{\boldsymbol{w}}^T \underline{\boldsymbol{x}}_n)$

• When
$$y_n = -1$$
, $\Pr(y_n = y'_n) = \operatorname{logit}(-\underline{\boldsymbol{w}}^T \underline{\boldsymbol{x}}_n)$



$$\operatorname{SCIHMS}(\operatorname{IPr}(y_n = y_n') = \operatorname{logit}(y_n \underline{\boldsymbol{w}}^T \underline{\boldsymbol{x}}_n)$$

Fitting Logistic Regression (Optional)

• Given:

$$\Pr(D|\underline{\boldsymbol{w}}) = \prod_{i=1}^{N} \operatorname{logit}(y_n \underline{\boldsymbol{w}}^T \underline{\boldsymbol{x}}_n)$$

- Re-write optimization as minimizing sum of point-wise errors $E_{in} = \sum e(h(\mathbf{x}_n), y_n)$
 - Maximize $\log(...)$ with respect to $\underline{\mathbf{w}}$ since \log argument > 0 and $\log()$ is monotone
- Equivalently, minimize:

$$\begin{split} -\frac{1}{N}\log(\ldots) &= -\frac{1}{N}\log(\prod(\ldots)) = -\frac{1}{N}\sum(\log(\ldots)) \\ &= \frac{1}{N}\sum\log(\frac{1}{\log\mathrm{it}(y_n\underline{\boldsymbol{w}}^T\underline{\boldsymbol{x}}_n)}) \\ &= \frac{1}{N}\sum\log(1+\exp(-y_n\underline{\boldsymbol{w}}^T\underline{\boldsymbol{x}}_n)) = \frac{1}{N}\sum e(h(\underline{\boldsymbol{x}}_n),y_n) = E_{in}(\underline{\boldsymbol{w}}) \end{split}$$



Gradient Descent for Logistic Regression

- Gradient descent requires two inputs:
 - Gradient of the cost function $\frac{\partial E}{w_i}$ for all j
 - Cost function $E_{in}(\underline{w})$
- The cost function is:

$$E_{in}(\underline{\boldsymbol{w}}) = \frac{1}{N} \sum_{i} e(h(\underline{\boldsymbol{x}}_{i}; \underline{\boldsymbol{w}}), y_{i})$$

The cost function for logistic regression:

$$E_{in}(\underline{\boldsymbol{w}}) = \frac{1}{N} \sum_{i} \log(1 + \exp(-y_i \underline{\boldsymbol{w}}^T \underline{\boldsymbol{x}}_i))$$

- Thus gradient descent converges to global minimum
 - It can be shown that $E_{in}(\underline{w})$ is convex in \underline{w}
 - In fact sum of exponentials and flipped exponentials is convex and log is

One-Vs-All Multi-Class Classification

- Aka "one-vs-rest" classifier
- Assume we have n classes $c_1, ..., c_n$ to distinguish given \underline{x}
- Learn
 - Create n binary classification problems where we classify c_i vs c_{-i} (everything but i)
 - Learn *n* classifiers with optimal $\underline{\mathbf{w}}_i$, each estimating $\Pr(y = i | \underline{\mathbf{x}}; \underline{\mathbf{w}}_i)$
- Predict
 - Evaluate the n classifiers
 - Pick the class y = i with the maximum $Pr(y = i | \underline{x})$



Cost Function for Multi-Class Probabilistic Classification

• The cost function for logistic regression is:

$$E_{in}(\underline{\boldsymbol{w}}) = -\frac{1}{N} \sum_{i=1}^{N} (y_i \log \Pr(y = 1 | \underline{\boldsymbol{x}}_i) + (1 - y_i) \log(1 - \Pr(y = 1 | \underline{\boldsymbol{x}}_i)))$$

- Encode the expected outputs y_i one-hot
 - I.e., the *j*-th element $\underline{\mathbf{y}}_{i}|_{j}$ is 1 iff the correct class is the *j*-th
 - E.g., for k = 4 1000
- Using $\underline{h}(\underline{x})$ as the outputs from the model and $\Pr(y=1|\underline{x})=p(\underline{x})$:

$$E_{in}(\underline{\boldsymbol{w}}) = -\frac{1}{N} \sum_{i=1}^{N} \sum_{k=1}^{K} \left(\underline{\boldsymbol{y}}_{i} \log(\underline{\boldsymbol{p}}(\underline{\boldsymbol{x}}_{i})) + (1 - \underline{\boldsymbol{y}}_{i}) \log(1 - \underline{\boldsymbol{p}}(\underline{\boldsymbol{x}}_{i})) \right) \Big|_{k}$$

- In the innermost summation we consider the error on each class / digit
 - E.g., for k = 4 1000 vs 0100



SCIENO the digits that are equal don't contribute to the error ACADE Thy different digits give a positive contribution

- Perceptron
- Logistic Regression
- LDA, QDA
- Kernel Methods



Basic Idea of Parametric Models

- We assume a model is generating the data
 - The functional form of the model is known
 - The model is parametrized with unknown parameters to estimate
- Pros
 - Utilize structure in the data
 - Models are easy to fit: few parameters to estimate
 - Accurate predictions if model assumptions are correct
- Cons
 - Strong assumptions about the data
 - Low accuracy if model assumptions are incorrect



Linear and Quadratic Discriminant Analysis

- Aka LDA and QDA
- Parametric models
 - Assume each class generating process is multivariate Gaussian
 - Classifiers with linear and quadratic decision surface
- Pros
 - Closed-form solutions easy to compute (sample mean and covariance)
 - Inherently multiclass
 - No hyperparams to tune
- Cons
 - Strong assumptions about the data



LDA / QDA: Model Form

- Both LDA and QDA assume that the class generating process $f_k(\underline{x}; \mu_k, \underline{\Sigma}_k)$ is from a multivariate Gaussian
- Linear discriminant analysis has a model:

$$f_k(\underline{x}; \underline{\mu_k}, \underline{\Sigma_k}) \sim \mathcal{N}(\underline{\mu_k}, \underline{\Sigma})$$

where:

- Means μ_k are different for all k classes
- The covariance matrix Σ is the same for all k classes
- It can be proven that the classes are separated by linear decision boundaries
- Quadratic discriminant analysis has a model:
 - Classes k have different covariance matrix Σ_k
 - It can be proved that the classes are separated by quadratic boundaries



Bayes Theorem for LDA / QDA

- ullet Consider a classification setup with multi-class output $Y \in \{1,...,K\}$
- We need to build a parametric model for the conditional distribution:

$$\Pr(Y = k | X = \underline{x})$$

• Use Bayes theorem:

$$\Pr(Y = k | X = \underline{x}) = \frac{\Pr(X = \underline{x} | Y = k) \cdot \Pr(Y = k)}{\Pr(X = \underline{x})}$$

where:

- $Pr(X = \underline{x} | Y = k)$: given a class, estimate the probability of \underline{x}
- $Pr(Y = k) = \pi_k$: estimate the probability of each class (prior)
- $Pr(X = \underline{x})$: estimate the probability of each input
- Estimate probabilities from data



LDA / QDA: Boundary Decision (*)

• Consider the ratio between the probabilities of Y = k vs Y = j:

$$r = \frac{\Pr(Y = k | X = \underline{x})}{\Pr(Y = j | X = \underline{x})}$$

• Using the model assumption and Bayes theorem:

$$\Pr(Y = k | X = \underline{x}) \propto \Pr(X = \underline{x} | Y = k) \cdot \Pr(Y = k)$$
$$= f_k(\underline{x}; \underline{\mu_k}) \cdot \pi_k$$

where:

$$f_k(\underline{\mathbf{x}}) = c \cdot \exp\left(-\frac{1}{2}(\underline{\mathbf{x}} - \underline{\mu_k})^T \underline{\underline{\boldsymbol{\Sigma}}}^{-1}(\underline{\mathbf{x}} - \underline{\mu_k})\right)$$

• We can apply a log(·) since it is a monotone transformation

$$r = \log \frac{f_k(\underline{\mathbf{x}})}{f_j(\underline{\mathbf{x}})} + \log \frac{\pi_k}{\pi_j}$$



LDA / QDA: Boundary Decision (*)

• We can apply a $log(\cdot)$ since it is a monotone transformation

$$r = \log \frac{f_k(\underline{\mathbf{x}})}{f_j(\underline{\mathbf{x}})} + \log \frac{\pi_k}{\pi_j}$$

- The second term does not depend on \underline{x}
- The first term is proportional to:

$$(\underline{\mathbf{x}} - \underline{\mu_k})^T \underline{\underline{\mathbf{\Sigma}}}^{-1} (\underline{\mathbf{x}} - \underline{\mu_k}) - (\underline{\mathbf{x}} - \underline{\mu_j})^T \underline{\underline{\mathbf{\Sigma}}}^{-1} (\underline{\mathbf{x}} - \underline{\mu_j})$$

• By expanding the expression, simplifying $\underline{x}^T \underline{\underline{\Sigma}}^{-1} \underline{x}$ noticing that $\underline{x}^T \underline{\underline{\Sigma}}^{-1} (\underline{\mu_k} - \underline{\mu_j})$ plus its transposed value (because $\underline{\underline{\Sigma}}$ is symmetrical) is equal to $2\underline{x}^T \underline{\underline{\Sigma}}^{-1} (\underline{\mu_k} - \underline{\mu_j})$ we get:

$$-\frac{1}{2}(\underline{\mu_k} + \underline{\mu_j})^T \underline{\underline{\Sigma}}^{-1}(\underline{\mu_k} - \underline{\mu_j}) + \underline{\mathbf{x}}^T \underline{\underline{\Sigma}}^{-1}(\underline{\mu_k} - \underline{\mu_j})$$



LDA / QDA: Learn Model

- In practice:
 - We don't care about $Pr(X = \underline{x})$ since this is common for all classes
 - We assume to know the prior $\Pr(Y = k) = \pi_k$ or we estimate it from the data
 - We need to estimate the conditional probability $Pr(X = \underline{x}|Y = k)$
- The model assumes that the conditional probability has a Gaussian distribution:

$$\Pr(X = \underline{x} | Y = k) = f_k(\underline{x}; \underline{\mu_k}, \underline{\Sigma_k})$$

where:

$$f_k(\underline{\boldsymbol{x}}) = \frac{1}{(2\pi)^n |\boldsymbol{\Sigma_k}|^{1/2}} \exp\left(-\frac{1}{2}(\underline{\boldsymbol{x}} - \underline{\mu_k})^T \underline{\underline{\boldsymbol{\Sigma}_k}}^{-1} (\underline{\boldsymbol{x}} - \underline{\mu_k})\right)$$

• We can estimate the parameters $\underline{\mu_k}$, $\underline{\underline{\Sigma_k}}$ from the data using sample mean and covariance



Evaluating LDA / QDA

• When we get a new $\underline{x} = \underline{x}'$ we compute for each class Y = k

$$\Pr(Y = k | X = \underline{x}) \propto f_k(\underline{x}; \mu_k, \underline{\Sigma_k}) \cdot \pi_k$$

and choose the k that maximizes the posterior probability

- Sometimes $f_i(\underline{x})$ is from a multivariate Gaussian distribution where $\underline{\underline{\Sigma}}$ is diagonal because of feature independence
 - Then the expressions can be simplified even more



- Perceptron
- Logistic Regression
- LDA, QDA
- Kernel Methods

