

# Selected Topics in Applied Machine Learning: An integrating view on data analysis and learning algorithms

ESLLI '2015  
Barcelona, Spain

<http://ufal.mff.cuni.cz/esslli2015>

Barbora Hladká  
hladka@ufal.mff.cuni.cz

Martin Holub  
holub@ufal.mff.cuni.cz

Charles University in Prague,  
Faculty of Mathematics and Physics,  
Institute of Formal and Applied Linguistics

# Block 3.1

## Ensemble learning methods

### Outline

- Combining classifiers into ensembles
- Bagging vs. boosting
- Bagging – example classifier
- Random Forests
- AdaBoost

# Ensemble classifiers – a motivation exercise

**Consider the following task** – we have a binary classification problem and a number of predictors, each with error less than 0.5. Will the resulting majority voting ensemble have an error lower than the single classifiers?

# Ensemble classifiers – a motivation exercise

**Consider the following task** – we have a binary classification problem and a number of predictors, each with error less than 0.5. Will the resulting majority voting ensemble have an error lower than the single classifiers?

– Depends on the *accuracy* and the *diversity* of the base learners!

# Ensemble classifiers – a motivation exercise

**Consider the following task** – we have a binary classification problem and a number of predictors, each with error less than 0.5. Will the resulting majority voting ensemble have an error lower than the single classifiers?

– Depends on the *accuracy* and the *diversity* of the base learners!

**Particular settings** – assume that you have

- 21 classifiers
- each with error  $p = 0.3$
- their outputs are *statistically independent*

**Compute the error of the ensemble under these conditions!**

## Resampling approach

- Distribute the training data into  $K$  portions
- Run the learning process to get  $K$  different models
- Collect the output of the  $K$  models use a combining function to get a final output value

# Bootstrapping principle

- New data sets  $Data_1, \dots, Data_K$  are drawn from  $Data$  with replacement, each of the same size as the original  $Data$ , i.e.  $n$ .
- In the  $i$ -th step of the iteration,  $Data_i$  is used as a training set, while the examples  $\{\mathbf{x} \mid \mathbf{x} \in Data \wedge \mathbf{x} \notin Data_i\}$  form the test set.

# Bootstrapping principle

- New data sets  $Data_1, \dots, Data_K$  are drawn from  $Data$  with replacement, each of the same size as the original  $Data$ , i.e.  $n$ .
- In the  $i$ -th step of the iteration,  $Data_i$  is used as a training set, while the examples  $\{\mathbf{x} \mid \mathbf{x} \in Data \wedge \mathbf{x} \notin Data_i\}$  form the test set.
- The probability that we pick an instance is  $1/n$ , and the probability that we do not pick an instance is  $1 - 1/n$ . The probability that we do not pick it after  $n$  draws is  $(1 - 1/n)^n \approx e^{-1} \doteq 0.368$ .
- It means that for training the system will not use 36.8% of the data, and the error estimate will be pessimistic. So the solution is to repeat the process many times.



# Same algorithm, different classifiers

Combining classifiers to improve the performance

## Bootstrapping methods – key ideas

- combining the classification results from different classifiers to produce the final output
- using (un)weighted voting
- different training data
- different features
- different values of the relevant parameters
- performance: complementarity  $\rightarrow$  potential improvement

## Two fundamental approaches

- **Bagging** works by taking a bootstrap sample from the training set
- **Boosting** works by changing the weights on the training set

# Bagging and boosting — the difference

- **Bagging:** each predictor is trained independently
- **Boosting:** each predictor is built on the top of previous predictors trained
  - Like bagging, boosting is also a voting method. In contrast to bagging, boosting actively tries to generate complementary learners by training the next learner on the mistakes of the previous learners.

# Combining multiple learners

- the more **complementary** the learners are, the more useful their combining is
- the simplest way to combine multiple learners is **voting**
- in **weighted voting** the voters (= base-learners) can have different weights

## Unstable learning

- learning algorithm is called unstable if small changes in the training set cause large differences in generated models
- typical unstable algorithm is the decision trees learning
- bagging or boosting techniques are a natural remedy for unstable algorithms

- Bagging is a voting method that uses slightly different training sets (generated by bootstrap) to make different base-learners. Generating complementary base-learners is left to chance and to instability of the learning method.
- Generally, bagging can be combined with any approach to learning.

## Bootstrap **AGG**regat**ING**

- 1 for  $i \leftarrow 1$  to  $K$  do
- 2  $Train_i \leftarrow \text{bootstrap}(Data)$
- 3  $h_i \leftarrow \text{TrainPredictor}(Train_i)$

## Combining function

- **Classification:**  $h_{final}(\mathbf{x}) = \text{MajorityVote}(h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_K(\mathbf{x}))$
- **Regression:**  $h_{final}(\mathbf{x}) = \text{Mean}(h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_K(\mathbf{x}))$

# Random Forests

- an ensemble method based on decision trees and bagging
- builds a number of random decision trees and then uses voting
- introduced by L. Breiman (2001), then developed by L. Breiman and A. Cutler
- very good (state-of-the-art) prediction performance
- a nice page with description  
[www.stat.berkeley.edu/~breiman/RandomForests/cc\\_home.htm](http://www.stat.berkeley.edu/~breiman/RandomForests/cc_home.htm)
- important: Random Forests helps to
  - avoid overfitting (by random sampling the training data set)
  - select important/useful features (by random sampling the feature set)

## The algorithm for building a tree in the ensemble

- 1 Having a training set of the size  $n$ , sample  $n$  cases at random - but with replacement, and use the sample to build a decision tree.
- 2 If there are  $M$  input features, choose a less number  $m \ll M$  (fixed for the whole procedure). When building the tree, at each node  $m$  variables are selected at random out of the  $M$  and the best split on these  $m$  features is used to split the node.
- 3 Each tree is grown to the largest extent possible. There is no pruning.

**The more trees in the ensemble, the better.  
There is no risk of overfitting!**

# Regularized Random Forests

- a recent extension of the original Random Forest
  - introduced by Houtao Deng and George Runger (2012)
- produces a compact feature subset
- provides an effective and efficient feature selection solution for many practical problems
- overcomes the weak spot of the ordinary RF: Random Forest importance score is biased toward the variables having more (categorical) values
- a useful page: <https://sites.google.com/site/houtaodeng/rrf>
  - a presentation
  - a sample code
  - links to papers
  - a brief explanation of the difference between RRF and guided RRF



# R packages for Random Forests

- **randomForest**: Breiman and Cutler's random forests for classification and regression
  - Classification and regression based on a forest of trees using random inputs.
- **RRF**: Regularized Random Forest
  - Feature Selection with Regularized Random Forest. This package is based on the 'randomForest' package by Andy Liaw. The key difference is the RRF function that builds a regularized random forest.
  - <http://cran.r-project.org/web/packages/RRF/index.html>
- **party**: A Laboratory for Recursive Partytioning
  - a computational toolbox for recursive partitioning
  - `cforest()` provides an implementation of Breiman's random forests
  - extensible functionality for visualizing tree-structured regression models is available

## Motivation

- I want to write a program that will accurately predict the winner of a tennis tournament based on the information like number of tournaments recently won by each player.
- I have not much experience so I ask a highly successful expert gambler to explain his betting strategy. In general, he is not able to explain a grand set of rules for predicting a winner. However, when he is provided with the data for a particular tournament, the expert has no problem to come up with a "rule of thumb" like *Bet on the player who has recently won the most matches*.
- Such a rule of thumb is obviously rough and inaccurate, we can expect to provide predictions that are at least a little bit better than random guessing.

## Motivation

- How to extract rules of thumb from expert that will be the most useful?
- How to combine moderately accurate rules of thumb into a single highly accurate prediction rule?

## Basic idea

- Boosting is a method that produces a very accurate predictor by combining rough and moderately accurate predictors.
- It is based on the observation that finding many rough predictors (rules of thumb) can be easier than finding a single, highly accurate predictor.

**AdaBoost** is a boosting method that repeatedly calls a given weak learner, each time with different distribution over the training data. Then we combine these weak learners into a strong learner.

- originally proposed by Freund and Schapire (1996)
- nice presentation including theoretical details and a demonstration available at  
[http://cmp.felk.cvut.cz/~sochmj1/adaboost\\_talk.pdf](http://cmp.felk.cvut.cz/~sochmj1/adaboost_talk.pdf)

## Key questions

- How to choose the distribution?
- How to combine the weak predictors into a single predictor?
- How many weak predictors should be trained?

**Schapire's strategy:** Change the distribution over the examples in each iteration, feed the resulting sample into the weak learner, and then combine the resulting hypotheses into a voting ensemble, which, in the end, would have a boosted prediction accuracy.

# Binary classification and AdaBoost

- We explain the notion of boosting using binary classification with the training data

$$Data = \{\langle \mathbf{x}_i, y_i \rangle : \mathbf{x}_i \in X, y_i \in Y, Y = \{-1, +1\}, i = 1, \dots, n\}$$

- We need to define distribution  $\mathcal{D}$  over  $Data$  such that  $\sum_{i=1}^n \mathcal{D}_i = 1$ .
- A weak classifier  $h_t : X \rightarrow Y$  has the property

$$\text{error}_{\mathcal{D}}(h_t) < 1/2.$$

# AdaBoost

- Initialize  $\mathcal{D}_1(i) = 1/n$
- At each step  $t$ 
  - Learn  $h_t$  using  $\mathcal{D}_t$ : find the weak classifier  $h_t$  with the minimum weighted sample error  $\text{error}_{\mathcal{D}_t}(h_t) = \sum_{i=1}^n \mathcal{D}_t(i) \delta(h(\mathbf{x}_i) \neq y_i)$
  - Set weight  $\alpha_t$  of  $h_t$  based on the sample error

$$\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \text{error}_{\mathcal{D}_t}(h_t)}{\text{error}_{\mathcal{D}_t}(h_t)} \right)$$

- Update the distribution ( $Z_t$  is a normalization factor)

$$\mathcal{D}_{t+1} = \frac{1}{Z_t} \mathcal{D}_t e^{-\alpha_t y_i h_t(\mathbf{x}_i)}$$

- Stop when impossible to find a weak classifier being better than chance
- Output the final classifier  $h_{final}(\mathbf{x}) = \text{sign} \sum_{i=1}^T \alpha_i h_i(\mathbf{x})$

- constructing  $\mathcal{D}_t$ :
  - On each round, the weights of incorrectly classified instances are increased so that the algorithm is forced to focus on the hard training examples.
  - $\mathcal{D}_1(i) = 1/n$
  - given  $\mathcal{D}_t$  and  $h_t$  (i.e. update  $\mathcal{D}_t$ ):

$$\mathcal{D}_{t+1}(i) = \frac{\mathcal{D}_t(i)}{Z_t} \cdot \begin{cases} e^{-\alpha_t} & \text{if } y_i = h_t(x_i) \\ e^{\alpha_t} & \text{if } y_i \neq h_t(x_i) \end{cases} = \frac{\mathcal{D}_t(i)}{Z_t} e^{-\alpha_t y_i h_t(x_i)},$$

where  $Z_t$  is normalization constant  $Z_t = \sum_i \mathcal{D}_t(i) e^{-\alpha_t y_i h_t(x_i)}$

- $\alpha_t$  measures the importance that is assigned to  $h_t$



## Weights

- $error_{\mathcal{D}_t}(h_t) < \frac{1}{2} \Rightarrow \alpha_t > 0$
- the smaller the error, the bigger the weight of the weak learner
- The bigger the weight, the more impact on the strong classifier:  
 $error_{\mathcal{D}_t}(h_1) < error_{\mathcal{D}_t}(h_2) \Rightarrow \alpha_1 > \alpha_2$

- $$\mathcal{D}_{t+1} = \frac{1}{Z_t} \mathcal{D}_t e^{-\alpha_t y_i h_t(\mathbf{x}_i)}$$

The weights of correctly classified instances are reduced while weights of misclassified instances are increased.

## Multiclass problem

- Assume classification task where  $Y = \{y_1, \dots, y_k\}$

$$h_t : X \rightarrow Y,$$

$$\mathcal{D}_{t+1}(i) = \frac{\mathcal{D}_t(i)}{Z_t} \cdot \begin{cases} e^{-\alpha_t} & \text{if } y_i = h_t(\mathbf{x}_i) \\ e^{\alpha_t} & \text{if } y_i \neq h_t(\mathbf{x}_i) \end{cases}$$

$$h_{final}(\mathbf{x}) = \operatorname{argmax}_{y \in Y} \sum_{\{y \mid h_t(\mathbf{x})=y\}} \alpha_t.$$

We can prove same bound on the error **if**  $\forall t : \epsilon_t \leq \frac{1}{2}$

# Block 3.2

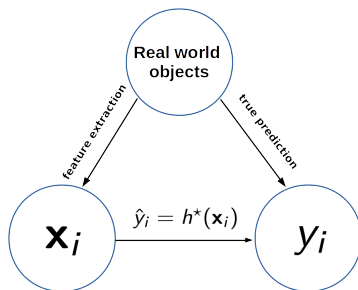
## Model complexity and regularization

### Outline

- Overfitting
- Regularization theoretically
  - Ridge regression
  - Lasso
  - Recap of linear regression
  - Recap of logistic regression

# Settings

- Suppose  $m$  features  $A_1, \dots, A_m$  and a set of possible target values  $Y$
- Suppose development data as a set of instances  $D = \{(\mathbf{x}_i, y_i), \mathbf{x}_i = \langle x_i^1, \dots, x_i^m \rangle, y_i \in Y\}$ , where  $\mathbf{x}_i$  are feature vectors and  $y_i$  are the corresponding target values

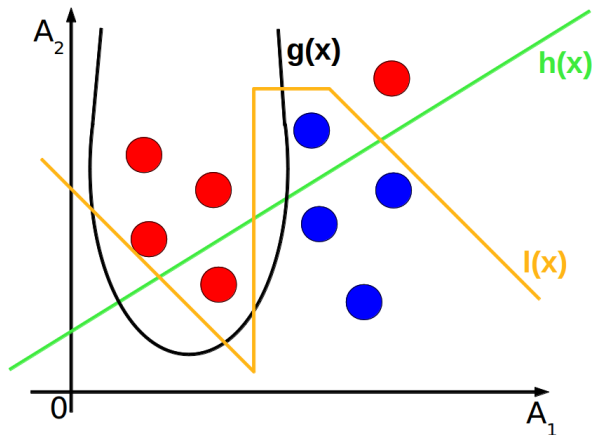


Let  $h^*$  be a best approximation of  $c$  trained on  $D$ .

**Model complexity** is the number of hypothesis parameters

$$\Theta = \langle \Theta_0, \dots, \Theta_m \rangle$$

# Model complexity – example

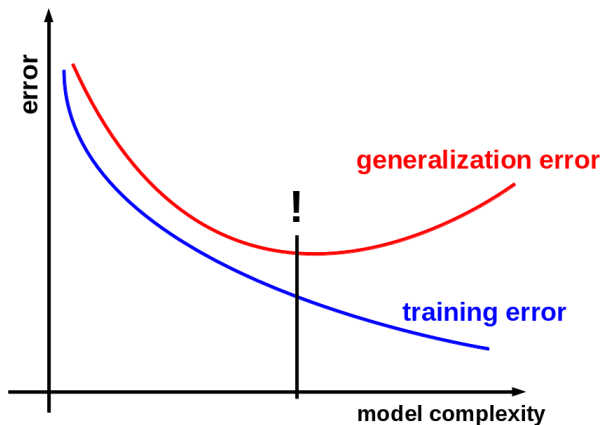


# Model complexity – example

- $h_1(\mathbf{x})$ : a straight line – determined by *two* parameters of the prediction function
  - doesn't fit two examples
- $h_2(\mathbf{x})$ : a parabola – determined by *three* parameters of the prediction function
  - doesn't fit one example
- $h_3(\mathbf{x})$ : a curve – determined by *many* parameters of the prediction function
  - perfectly fits all examples

# Model complexity and overfitting

Finding a model that minimizes generalization error  
... is one of central goals of the machine learning process

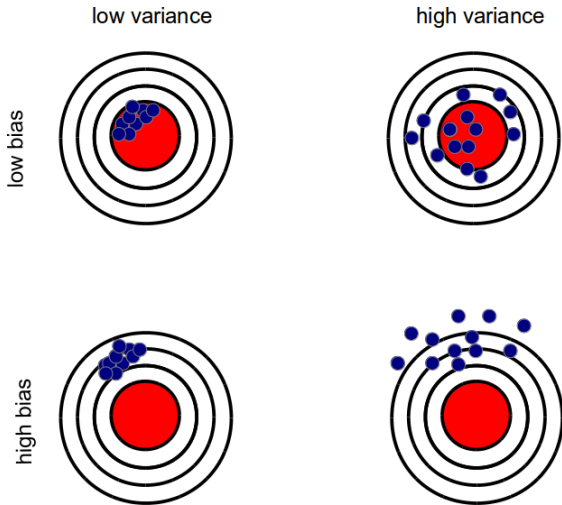




# Bias and variance

- 1 Select a machine learning algorithm
  - 2 Get  $k$  different training sets
  - 3 Get  $k$  predictors  $h_1^*, \dots, h_k^*$
- **Bias** measures error that originates from the learning algorithm  
– how far off in general the predictions by  $k$  predictors are from the true output value
  - **Variance** measures error that originates from the training data  
– how much the predictions for a test instance vary between  $k$  predictors

# Bias and variance



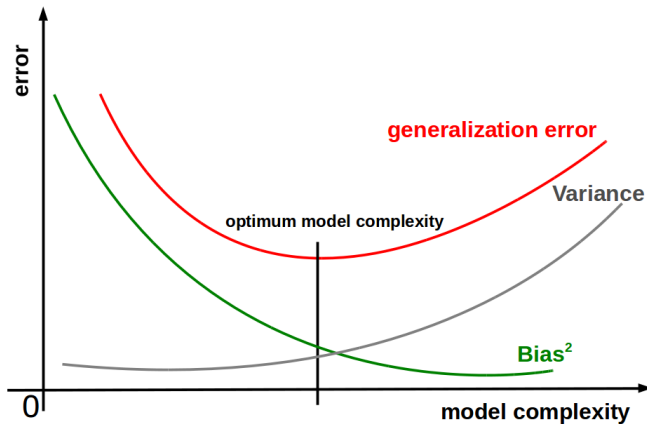
Generalization error  $error_{\mathcal{D}}(h)$  measures how well a hypothesis  $h$  generalizes beyond the used training data set, to unseen data with distribution  $\mathcal{D}$ .

**Decomposition of  $error_{\mathcal{D}}(h)$**

$$error_{\mathcal{D}}(h) = \text{Bias}^2 + \text{Variance}$$

# Bias and variance

- underfitting = high bias
- overfitting = high variance



We want a model in between which is

- powerful enough to model the underlying structure of data
- not so powerful to model the structure of the training data

Let's prevent overfitting by **complexity regularization**, a technique that regularizes the parameter estimates, or equivalently, shrinks the parameter estimates towards zero

- A machine learning algorithm estimates hypothesis parameters  $\Theta = \langle \Theta_0, \Theta_1, \dots, \Theta_m \rangle$  using  $\Theta^*$  that minimizes **loss** function for the data  $D$

$$\Theta^* = \underset{\Theta}{\operatorname{argmin}} \operatorname{loss}(\Theta)$$

- **Regularization**

$$\Theta^* = \underset{\Theta}{\operatorname{argmin}} \operatorname{loss}(\Theta) + \lambda * \mathbf{penalty}(\Theta)$$

where  $\lambda \geq 0$  is a **tuning parameter**

# Regularization – Ridge regression

$$\text{penalty}(\Theta) = \Theta_1^2 + \dots + \Theta_m^2$$

$$\Theta^* = \underset{\Theta}{\operatorname{argmin}} \operatorname{loss}(\Theta) + \lambda * (\Theta_1^2 + \dots + \Theta_m^2)$$

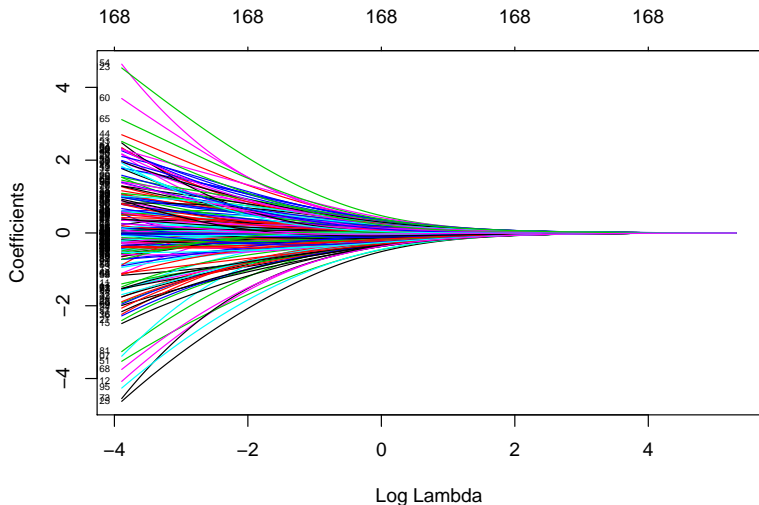
The penalty is applied to  $\Theta_1, \dots, \Theta_m$ , but not to  $\Theta_0$ , since the goal is to regularize the estimated association between each feature and the target value.

$$\Theta^* = \underset{\Theta}{\operatorname{argmin}} \operatorname{loss}(\Theta) + \lambda * (\Theta_1^2 + \dots + \Theta_m^2)$$

- **Let**  $\Theta_{\lambda_1}^*, \dots, \Theta_{\lambda_m}^*$  be ridge regression parameter estimates for a particular value of  $\lambda$
- **Let**  $\Theta_1^*, \dots, \Theta_m^*$  be unregularized parameter estimates
- **When**  $\lambda = 0$ , **then**  $\Theta_{\lambda_i}^* = \Theta_i^*$  for  $i = 1, \dots, m$
- **When**  $\lambda$  is extremely large, **then**  $\Theta_{\lambda_i}^* = 0$  for  $i = 1, \dots, m$
- **When**  $\lambda$  between, we are fitting a model and shrinking the parameters



# Ridge regression

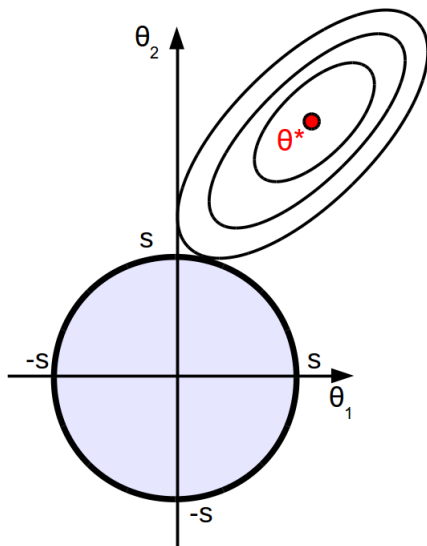


# Ridge regression – alternative formulation

$$\Theta^* = \operatorname{argmin}_{\Theta} \operatorname{loss}(\Theta)$$

subject to  $\Theta_1^2 + \dots + \Theta_m^2 \leq s$

- the gray circle represents the feasible region for Ridge regression; the contours represent different loss values for the unconstrained model

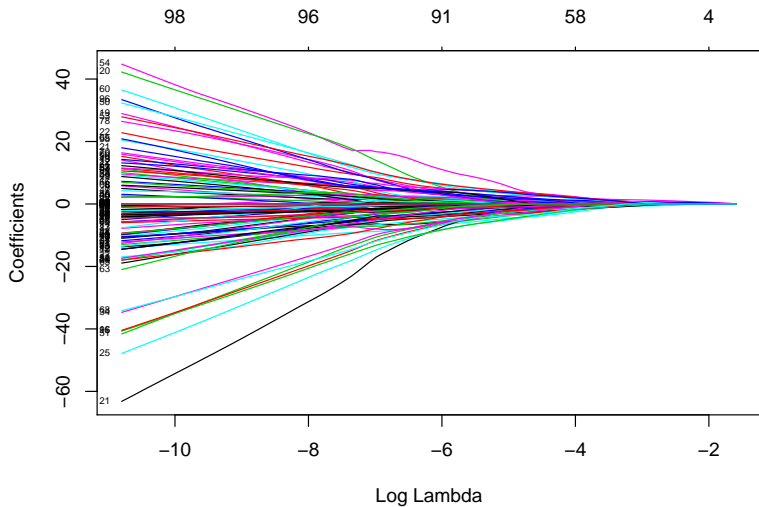


$$\text{penalty}(\Theta) = |\Theta_1| + \dots + |\Theta_m|$$

$$\Theta^* = \underset{\Theta}{\text{argmin}} \text{loss}(\Theta) + \lambda * (|\Theta_1| + \dots + |\Theta_m|)$$

$$\Theta^* = \underset{\Theta}{\operatorname{argmin}} \operatorname{loss}(\Theta) + \lambda * (|\Theta_1| + \dots + |\Theta_m|)$$

- **Let**  $\Theta_{\lambda 1}^*, \dots, \Theta_{\lambda m}^*$  be lasso regression parameter estimates
- **Let**  $\Theta_1^*, \dots, \Theta_m^*$  be unregularized parameter estimates
- **When**  $\lambda = 0$ , **then**  $\Theta_{\lambda i}^* = \Theta_i^*$  for  $i = 1, \dots, m$
- **When**  $\lambda$  grows, **then** the impact of penalty grows
- **When**  $\lambda$  is extremely large, **then**  $\Theta_{\lambda i}^* = 0$  for  $i = 1, \dots, m$

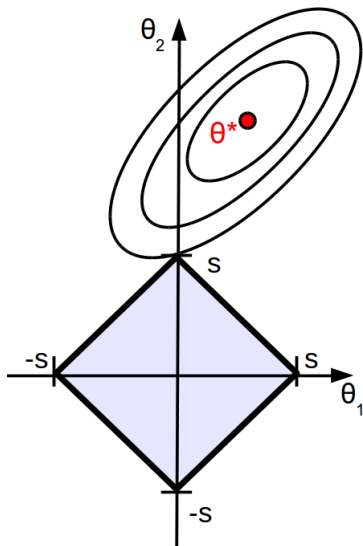


# Lasso – alternative formulation

$$\Theta^* = \underset{\Theta}{\operatorname{argmin}} \operatorname{loss}(\Theta)$$

subject to  $|\Theta_1| + \dots + |\Theta_m| \leq s$

- the grey square represents the feasible region of the Lasso; the contours represent different loss values for the unconstrained model
- the feasible point that minimizes the loss is more likely to happen on the coordinates on the Lasso graph than on the Ridge regression graph since the Lasso graph is more angular



## Difference between L1 and L2

Ridge regression shrinks all the parameters but eliminates none, while the Lasso can shrink some parameters to zero.

# Recap of linear regression

**Linear regression** is a regression algorithm

$$\Theta^* = \underset{\Theta}{\operatorname{argmin}} \sum_{i=1}^n (h(\mathbf{x}_i) - y_i)^2$$

where

- $h(\mathbf{x}) = \Theta_0 + \Theta_1 x_1 + \dots + \Theta_m x_m$
- loss function = mean squared error



## Intepretation of $\Theta$

- $h(\mathbf{x}) = \Theta_0 + \Theta_1 x_1 + \dots + \Theta_m x_m$

$\Theta_j$  gives an average change in a target value with one-unit change in feature  $A_j$ , holding other features fixed

# Regularized linear regression

$$h(\mathbf{x}) = \Theta_0 + \Theta_1 x_1 + \cdots + \Theta_m x_m$$

$$\Theta^* = \underset{\Theta}{\operatorname{argmin}} \sum_{i=1}^n (h(\mathbf{x}_i) - y_i)^2 + \lambda * \operatorname{penalty}(\Theta)$$

# Recap of logistic regression

**Logistic regression** is a classification algorithm

Assume  $Y = \{0, 1\}$

- **modeling the probability**  $h(\mathbf{x}) = \Pr(Y = 1|\mathbf{x}; \Theta)$

$$h(\mathbf{x}) = g(\Theta^T \mathbf{x}) = \frac{1}{1 + e^{-\Theta^T \mathbf{x}}}, \text{ where } \Theta = \langle \Theta_0, \dots, \Theta_m \rangle$$

- **prediction function** of  $\mathbf{x}$

$$= \begin{cases} 0, & h(\mathbf{x}) \geq 0.5 \\ 1, & h(\mathbf{x}) < 0.5 \end{cases}$$

# Recap of logistic regression

- $\frac{h(\mathbf{x})}{1 - h(\mathbf{x})} = \text{odds ratio}$
- **log odds is linear**

$$\log \frac{h(\mathbf{x})}{1 - h(\mathbf{x})} = \Theta^T \mathbf{x}$$

- **recall linear regression**

$$h(\mathbf{x}) = \Theta^T \mathbf{x}$$

## Interpretation of $\Theta$

Suppose  $\Theta = \langle \Theta_0, \Theta_1 \rangle$

- linear regression  $h(\mathbf{x}) = \Theta_0 + \Theta_1 x_1$ :  $\Theta_1$  gives an average change in a target value with one-unit change in  $A_1$
- logistic regression  $\log \frac{h(\mathbf{x})}{1-h(\mathbf{x})} = \Theta_0 + \Theta_1 x_1$ :  $\Theta_1$  gives an average change in  $\logit h(\mathbf{x})$  with one-unit change in  $A_1$

# Recap of logistic regression

## Interpretation of $\Theta$

**Example:** Classify CRY into two classes: "1"  $\sim$  "1", "0"  $\sim$  others. Use one feature only. Thus

$$\log \frac{h(\mathbf{x})}{1 - h(\mathbf{x})} = \Theta_0 + \Theta_1 x_1$$

Let  $p_1 = \Pr(Y = 1 | x_1 = 0)$  and  $p_2 = \Pr(Y = 1 | x_1 = 1)$ . Then

- $\log \frac{p_1}{1-p_1} = \Theta_0 + \Theta_1 x_1 \xrightarrow{x_1=0} \frac{p_1}{1-p_1} = e^{\Theta_0}$
- $\log \frac{p_2}{1-p_2} = \Theta_0 + \Theta_1 x_1 \xrightarrow{x_1=1} \frac{p_2}{1-p_2} = e^{\Theta_0 + \Theta_1}$
- $\frac{p_2}{1-p_2} / \frac{p_1}{1-p_1} = e^{\Theta_1}$  (i.e. the change in the odds of  $h(\mathbf{x})$  by unit change in  $x_1$ )

# Recap of logistic regression

## Estimating $\Theta$ by maximizing the likelihood

- likelihood of the data

$$\mathcal{L}(y_1, \dots, y_n; \Theta, \mathbf{X}) = \prod_{i=1}^n P(y_i | \mathbf{x}_i; \Theta)$$

- log likelihood of the data

$$\begin{aligned}\ell(y_1, \dots, y_n; \Theta, \mathbf{X}) &= \log L(y_1, \dots, y_n; \Theta, \mathbf{X}) \\ &= \sum_{i=1}^n \log P(y_i | \mathbf{x}_i; \Theta) \\ &= \sum_{i=1}^n y_i \log P(y_i | \mathbf{x}_i; \Theta) + (1 - y_i) \log(1 - P(y_i | \mathbf{x}_i; \Theta))\end{aligned}$$

# Recap of logistic regression

## Estimating $\Theta$ by maximizing the likelihood

- loss function

$$\begin{aligned} J(\Theta) &= \ell(y_1, \dots, y_n; \Theta, \mathbf{X}) \\ &= \sum_{i=1}^n y_i \log P(y_i | \mathbf{x}_i; \Theta) + (1 - y_i) \log(1 - P(y_i | \mathbf{x}_i; \Theta)) \end{aligned}$$

- optimization task

$$\begin{aligned} \Theta^* &= \operatorname{argmax}_{\Theta} J(\Theta) \\ &= \operatorname{argmin}_{\Theta} -J(\Theta) \\ &= \operatorname{argmin}_{\Theta} \sum_{i=1}^n -y_i \log P(y_i | \mathbf{x}_i; \Theta) - (1 - y_i) \log(1 - P(y_i | \mathbf{x}_i; \Theta)) \end{aligned}$$



**Multinomial logistic regression**  $Y = \{y_1, \dots, y_k\}$

- train  $k$  one-versus-all binary classifiers  $h_i^*, i = 1, \dots, k$
- classify  $\mathbf{x}$  into the class  $K$  that maximizes  $h_K^*(\mathbf{x})$

# Regularized logistic regression

$$h(\mathbf{x}) = \frac{1}{1 + e^{-\Theta^T \mathbf{x}}}$$

$$\Theta^* = \underset{\Theta}{\operatorname{argmin}} -\ell(y_1, \dots, y_n; \Theta, \mathbf{X}) + \lambda * \operatorname{penalty}(\Theta)$$

We will address the MOV task using

- 1 linear regression
- 2 regularized linear regression

We will address the VPR task using

- 1 logistic regression
- 2 regularized logistic regression