

Data Analysis and Knowledge Discovery

Linear models

Antti Airola

University of Turku
Department of Computing

antti.airola@utu.fi

- ▶ Given: *training set* of input-output pairs
 $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- ▶ Learn: model f such that given a new input \mathbf{x} for which corresponding output y unknown
 - ▶ $f(\mathbf{x}) \approx y$
- ▶ $\mathbf{x} \in \mathbb{R}^d$ is a d -dimensional feature vector (input)
- ▶ $y \in \mathbb{R}$ is the real-valued output to be predicted

The linear model, a simple yet popular choice:

$$f(\mathbf{x}) = w_1 \cdot x_1 + \dots + w_d \cdot x_d + b$$

- ▶ x_1, \dots, x_d , feature values
- ▶ w_1, \dots, w_d model coefficients
- ▶ $b \in \mathbb{R}$ intercept term

In sum notation, this is written as

$$f(\mathbf{x}) = \sum_{i=1}^d w_i \cdot x_i + b$$

- ▶ x_1, \dots, x_d , feature values
- ▶ w_1, \dots, w_d model coefficients
- ▶ $b \in \mathbb{R}$ intercept term

Append value 1 to the beginning of each feature vector (new constant valued feature x_0), and define a new coefficient $w_0 = b$

$$f(\mathbf{x}) = \sum_{i=1}^d w_i \cdot x_i + b = \sum_{i=0}^d w_i \cdot x_i$$

- ▶ $x_0 = 1$ the constant valued feature
- ▶ x_1, \dots, x_d , feature values
- ▶ w_0, \dots, w_d model coefficients

(A standard trick, re-naming the bias term 'b' just makes the following math and algorithmics a bit simpler.)

Finally, defining a coefficient vector $\mathbf{w} = [w_0, \dots, w_d]$, we can reformulate this as the inner product between the model and feature vectors

$$f(\mathbf{x}) = \sum_{i=0}^d w_i \cdot x_i = \mathbf{w}^T \mathbf{x}$$

- ▶ \mathbf{x} feature vector
- ▶ \mathbf{w} model vector

Regression line (single feature case)

given: A data set for two continuous attributes x (input feature) and y (output).

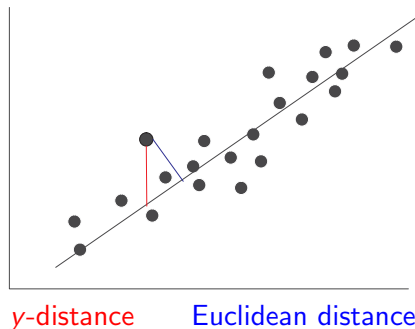
It is assumed that there is an approximate linear dependency between x and y :

$$y \approx w_1 x + w_0$$

Find a **regression line** (i.e. determine the parameters w_1 and w_0) such that the line fits the data as good as possible.

What is a **good fit**?

Regression line (single feature case)



Usually, the **mean square error in y-direction** is chosen as error measure (to be minimized).

It is equivalent to minimize the sum of squared errors in y-direction.

Given data (\mathbf{x}_i, y_i) ($i = 1, \dots, n$), the least squares error function is

$$\sum_{i=1}^n \left(\mathbf{w}^T \mathbf{x}_i - y_i \right)^2$$

Classical least-squares method (1809, Carl Friedrich Gauss):

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n \left(\mathbf{w}^T \mathbf{x}_i - y_i \right)^2$$

- ▶ choose model \mathbf{w}^* having smallest least-squares error
- ▶ can sometimes work fine if dimensionality d much smaller than training set size n
- ▶ prone to overfitting in high dimensions (also, no unique solution if $d > n$)
- ▶ Sensitive to outliers

Ridge regression, aka regularized least-squares:

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \left(\mathbf{w}^T \mathbf{x}_i - y_i \right)^2 + \lambda \sum_{i=0}^d w_i^2 \right\}$$

- ▶ regularization term penalizes too complex models
- ▶ $\lambda > 0$ regularization parameter (can be chosen with cross-validation)
- ▶ unique solution, much more robust than basic least-squares fitting, especially for high-dimensional data

Ridge regression:

$$\operatorname{argmin}_{\mathbf{w}} \left\{ \underbrace{\sum_{i=1}^n \left(\mathbf{w}^T \mathbf{x}_i - y_i \right)^2}_{\text{Training set error}} + \lambda \underbrace{\sum_{i=0}^d w_i^2}_{\text{Regularizer}} \right\}$$

Regularization for linear regression

Let us assume the following data structures:

- ▶ \mathbf{X} : a $n \times d$ -sized data matrix, one row for each instance
- ▶ \mathbf{y} : a n -length column vector of correct outputs, one element for each instance
- ▶ \mathbf{w} : d -length column vector of coefficients, we wish to learn from data
- ▶ (Minor technical detail: if we use the intercept term in our model, \mathbf{X} has one additional column of ones, and \mathbf{w} correspondingly one more element. This can be useful especially on low-dimensional data, where the additional flexibility given to model may be helpful.)

Regularization for linear regression

Ridge regression:

$$\operatorname{argmin}_{\mathbf{w}} \left\{ \underbrace{\sum_{i=1}^n \left(\mathbf{w}^T \mathbf{x}_i - y_i \right)^2}_{\text{Training set error}} + \lambda \underbrace{\sum_{i=0}^d w_i^2}_{\text{Regularizer}} \right\}$$

The same in matrix form:

$$\operatorname{argmin}_{\mathbf{w}} \left\{ \underbrace{(\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})}_{\text{Training set error}} + \underbrace{\lambda \mathbf{w}^T \mathbf{w}}_{\text{Regularizer}} \right\}$$

Solving ridge regression

Objective function to be minimized

$$J(\mathbf{w}) = (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda\mathbf{w}^T\mathbf{w}$$

Gradient

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial}{\partial \mathbf{w}} \left((\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda\mathbf{w}^T\mathbf{w} \right)$$

It can be shown that the global minimum of the objective function can be found at the point where the gradient is zero (due to convexity of the objective function).

Solving ridge regression

$$\begin{aligned}\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \left((\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w} \right) \\ &= \frac{\partial}{\partial \mathbf{w}} \left(\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y} + \lambda \mathbf{w}^T \mathbf{w} \right)\end{aligned}$$

Recall: rules of matrix transposition

$$(\mathbf{MN})^T = \mathbf{N}^T \mathbf{M}^T$$

Solving ridge regression

$$\begin{aligned}\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \left((\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w} \right) \\ &= \frac{\partial}{\partial \mathbf{w}} \left(\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y} + \lambda \mathbf{w}^T \mathbf{w} \right) \\ &= \frac{\partial}{\partial \mathbf{w}} \left(\mathbf{w}^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y} \right)\end{aligned}$$

Solving ridge regression

$$\begin{aligned}\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \left((\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^T \mathbf{w} \right) \\ &= \frac{\partial}{\partial \mathbf{w}} \left(\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y} + \lambda \mathbf{w}^T \mathbf{w} \right) \\ &= \frac{\partial}{\partial \mathbf{w}} \left(\mathbf{w}^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y} \right) \\ &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} + (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^T \mathbf{w} - 2\mathbf{X}^T \mathbf{y}\end{aligned}$$

Gradient rules

$$\frac{\partial}{\partial \mathbf{w}} \mathbf{v}^T \mathbf{w} = \mathbf{v}$$

$$\frac{\partial}{\partial \mathbf{w}} \mathbf{w}^T \mathbf{M} \mathbf{w} = \mathbf{M} \mathbf{w} + \mathbf{M}^T \mathbf{w}$$

Solving ridge regression

$$\begin{aligned}\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \left((\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^\top \mathbf{w} \right) \\ &= \frac{\partial}{\partial \mathbf{w}} \left(\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \mathbf{w} + \mathbf{y}^\top \mathbf{y} + \lambda \mathbf{w}^\top \mathbf{w} \right) \\ &= \frac{\partial}{\partial \mathbf{w}} \left(\mathbf{w}^\top (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} - 2\mathbf{y}^\top \mathbf{X} \mathbf{w} + \mathbf{y}^\top \mathbf{y} \right) \\ &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} + (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^\top \mathbf{w} - 2\mathbf{X}^\top \mathbf{y} \\ &= 2(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} - 2\mathbf{X}^\top \mathbf{y}\end{aligned}$$

Solving ridge regression

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 2(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})\mathbf{w} - 2\mathbf{X}^T \mathbf{y} = 0$$

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})\mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

Ridge regression

- ▶ Given: data matrix \mathbf{X} , outputs \mathbf{y} , hyperparameter λ
- ▶ Solve the above linear system to find model coefficients \mathbf{w}
- ▶ Can be solved using any standard linear algebra package (e.g. `numpy.linalg.solve()`)
- ▶ Computational complexity: $O(d^3 + d^2 n)$, memory usage $O(d^2 + dn)$
- ▶ feasible if dimensionality d is not too large (at most couple of thousands), but what if $d \gg n$?

Code for training ridge regression

```
import numpy as np

def ridge(X, y, regparam):
    #X: nxd data matrix
    #y: n-length vector of outputs
    #regparam > 0: parameter
    #returns: coefficients w
    d = X.shape[1]
    I = np.eye(d)
    A = X.T@X + regparam*I
    b = X.T@y
    w = np.linalg.solve(A, b)
    return w
```

Solving ridge regression, the other way

$$\begin{aligned}(\mathbf{X}\mathbf{X}^T + \lambda\mathbf{I})\mathbf{a} &= \mathbf{y} \\ \mathbf{w} &= \mathbf{X}^T\mathbf{a}\end{aligned}$$

Ridge regression

- ▶ It can be shown, that ridge regression can equivalently be trained by solving a $n \times n$ -sized linear system of equations
- ▶ Much more efficient than previous form, if dimensionality much larger than sample size (e.g. microarray data in bioinformatics)
- ▶ Computational complexity: $O(n^3 + n^2d)$, memory usage $O(n^2 + dn)$

Code for ridge regression, dual form

```
import numpy as np

def ridge(X, y, regparam):
    #X: nxd data matrix
    #y: n-length vector of outputs
    #regparam > 0: parameter
    #returns: coefficients w
    n = X.shape[0]
    I = np.eye(n)
    A = X@X.T + regparam*I
    a = np.linalg.solve(A, y)
    w = X.T@a
    return w
```


Code for lazy people

```
from sklearn.linear_model import Ridge

def rls_using_sklearn(X, y, regparam):
    learner = Ridge(alpha=regparam,
                     fit_intercept=False)
    learner.fit(X, y)
    w = learner.coef_
    return w
```

Should give same results as previous codes, the underlying implementation automatically decides, which of the previously shown formulations is solved depending on the values of d and n . Minor technical detail: Usually you should set `fit_intercept=True`, this would be equivalent to appending a constant feature to each feature vector in the previous examples.

Applying the model on new data

```
import numpy as np

def predict(x_test, w):
    #x_test: test instance, vector of d-features
    #w: vector of d-coefficients
    return x_test@w
```

What did we learn?

- ▶ linear model for regression
- ▶ ridge regression, aka regularized least squares, aka least-squares support vector machine
- ▶ idea: minimize mean squared error on training data, use a regularization term to penalize model complexity (here coefficients squared)
- ▶ unique optimal solution by solving a linear system of equations
- ▶ efficient to train, either by solving $d \times d$ or $n \times n$ -sized linear system (choose minimum)
- ▶ produces a compact linear model, that can be very efficiently be used to predict on test instances
- ▶ often gives good predictive performance, assuming the problem is not highly non-linear (you could also try, say, k-nearest neighbour, in case it is)

About model selection

- ▶ challenge: need to do model selection (default choices like $\lambda = 1$ may often work suboptimally or not at all)
- ▶ regularization parameter λ , value needs to be chosen for example by 10-fold or leave-one-out cross-validation
 - ▶ Advanced topic: very fast cross-validation algorithms exist for ridge regression (especially fast leave-one-out widely known and found in most decent implementations)
 - ▶ My rule of thumb: in most cases selecting λ by choosing the parameter leading to lowest cross-validation error from the exponential grid $\{2^{-15}, \dots, 2^{15}\}$ should suffice.

Wait, what if we had chosen to minimize some other reasonable criterion?

- ▶ for example, what if we would measure the fit of model to training data with absolute error $|f(\mathbf{x}) - y|$ instead of squared error $(f(\mathbf{x}) - y)^2$?
- ▶ or maybe I could penalize the model coefficients instead of the squared terms w_i^2 with, say, absolute magnitudes $|w_i|$?
- ▶ or maybe...

Ridge regression:

$$\operatorname{argmin}_{\mathbf{w}} \left\{ \underbrace{\sum_{i=1}^n \left(\mathbf{w}^T \mathbf{x}_i - y_i \right)^2}_{\text{Training set error}} + \lambda \underbrace{\sum_{i=0}^d w_i^2}_{\text{Regularizer}} \right\}$$

Lasso:

$$\operatorname{argmin}_{\mathbf{w}} \left\{ \underbrace{\sum_{i=1}^n \left(\mathbf{w}^T \mathbf{x}_i - y_i \right)^2}_{\text{Training set error}} + \lambda \underbrace{\sum_{i=0}^d |w_i|}_{\text{Regularizer}} \right\}$$

- ▶ Least absolute shrinkage and selection operator
- ▶ Lasso is a popular embedded feature selection algorithm
- ▶ when λ is large, many coefficients w_i tend to become zero
- ▶ ridge regression does not usually lead to zero coefficients!

Ridge regression:

$$\operatorname{argmin}_{\mathbf{w}} \left\{ \underbrace{\sum_{i=1}^n \left(\mathbf{w}^T \mathbf{x}_i - y_i \right)^2}_{\text{Training set error}} + \underbrace{\lambda \sum_{i=0}^d w_i^2}_{\text{Regularizer}} \right\}$$

Elastic Net:

$$\operatorname{argmin}_{\mathbf{w}} \left\{ \underbrace{\sum_{i=1}^n \left(\mathbf{w}^T \mathbf{x}_i - y_i \right)^2}_{\text{Training set error}} + \underbrace{\lambda_1 \sum_{i=0}^d |w_i| + \lambda_2 \sum_{i=0}^d w_i^2}_{\text{Regularizer}} \right\}$$

Ridge regression:

$$\operatorname{argmin}_{\mathbf{w}} \left\{ \underbrace{\sum_{i=1}^n \left(\mathbf{w}^T \mathbf{x}_i - y_i \right)^2}_{\text{Training set error}} + \lambda \underbrace{\sum_{i=0}^d w_i^2}_{\text{Regularizer}} \right\}$$

Support vector regression

$$\operatorname{argmin}_{\mathbf{w}} \left\{ \underbrace{\sum_{i=1}^n \max \left(0, |\mathbf{w}^T \mathbf{x}_i - y_i| - \epsilon \right)}_{\text{Training set error}} + \lambda \underbrace{\sum_{i=0}^d w_i^2}_{\text{Regularizer}} \right\}$$

- ▶ Why you might need to consider these alternatives
 - ▶ Lasso: leads to sparse solutions with many zero-coefficients, good for feature selection (`sklearn.linear_model.Lasso`)
 - ▶ Elastic Net: tries to get best of both worlds by interpolating between basic ridge regression and Lasso (`sklearn.linear_model.ElasticNet`)
 - ▶ Support vector regression: more robust towards outliers in data (`sklearn.svm.SVR`)
- ▶ Then again...
 - ▶ much more difficult to optimize, there is no analytical solution for minimizer, and since these are non-smooth basic gradient descent methods will not work well
 - ▶ in practice, often will not yield much better predictive accuracy
- ▶ Huge number of other approaches also out there

Classification as regression

- ▶ Given: *training set* of input-output pairs
 $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- ▶ Learn: model f such that given a new input \mathbf{x} for which corresponding output y unknown
 - ▶ $f(\mathbf{x}) \approx y$
- ▶ $\mathbf{x} \in \mathbb{R}^d$ is a d -dimensional feature vector (input)
- ▶ $y \in -1, 1$ denotes, whether instance belongs to class -1 or 1

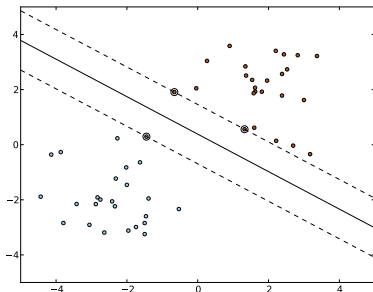
Classification as regression

A two-class classification problem (with classes encoded as -1 and 1) can be viewed as regression problem.

The regression function will usually not yield exact outputs -1 and 1, but the classification decision can be made by considering 0 as a cut-off value.

Problem: The objective functions aims at minimizing the function approximation error (for example, the mean squared error), but not misclassifications.

Linear classifier



- ▶ linear model defines a hyperplane $\{\mathbf{x} | \mathbf{w}^T \mathbf{x} + b = 0\}$
- ▶ this is the *decision boundary* between the classes
- ▶ $\{\mathbf{x} | \mathbf{w}^T \mathbf{x} + b > 0\}$ classified as positive
- ▶ $\{\mathbf{x} | \mathbf{w}^T \mathbf{x} + b < 0\}$ classified as negative
- ▶ may be fitted with ridge regression, but least squares loss does not optimize separability of the classes directly
 - ▶ sensitive to outliers - still, often works well in practice

- ▶ ridge regression is based on optimizing least squares fit of the model
- ▶ not optimal choice: in classification we only care if prediction falls on the correct side of the decision boundary!
- ▶ if true class $+1$: then predictions 0.1 , 1 , 10 , 10^6 all give correct classification
- ▶ least-squares loss tries to force all predictions to $+1/-1$
- ▶ can penalize arbitrarily much correct predictions with too large magnitude!

Loss functions

A loss function $L(\hat{y}, y)$ measures the price we pay for predicting \hat{y} when the true label is y

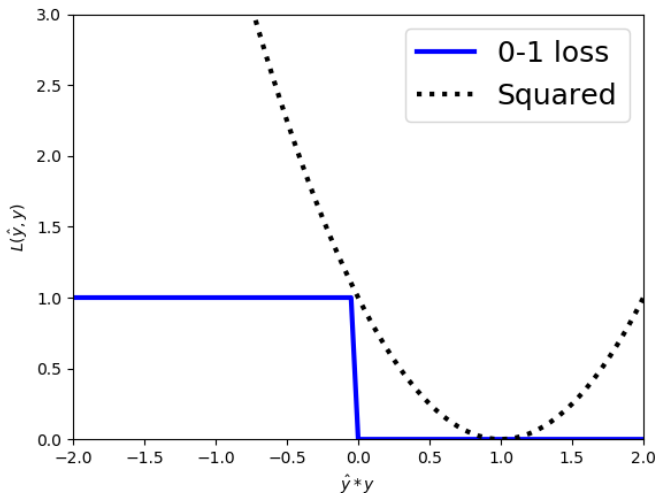
- ▶ *0-1 loss*: natural choice for classification

$$L(\hat{y}, y) = \begin{cases} 1, & \text{if } y * \hat{y} < 0 \\ 0, & \text{else} \end{cases}$$

- ▶ *squared loss*: used by ridge regression

$$L(\hat{y}, y) = (\hat{y} - y)^2$$

Squared loss



Support vector machine

- ▶ can we do better?
- ▶ optimizing 0-1 loss directly would lead to intractable optimization problem
- ▶ reasonable convex approximation: the hinge loss
- ▶ leads to the method of *support vector machine*
- ▶ geometric interpretation: maximises the margin separating the classes

Loss functions

A loss function $L(\hat{y}, y)$ measures the price we pay for predicting \hat{y} when the true label is y

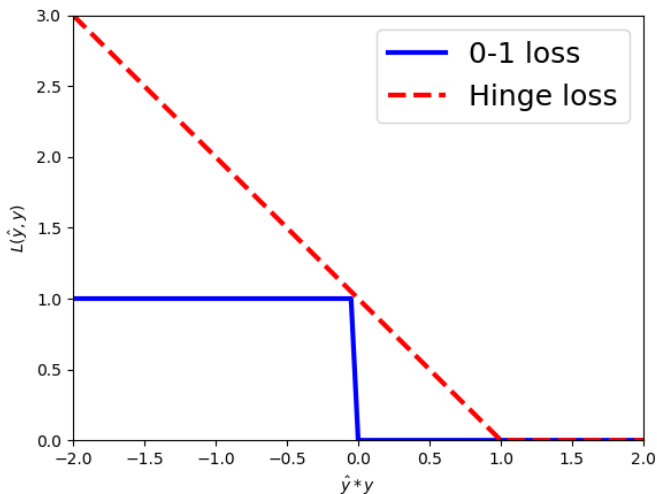
- ▶ *0-1 loss*: natural choice for classification

$$L(\hat{y}, y) = \begin{cases} 1, & \text{if } y * \hat{y} < 0 \\ 0, & \text{else} \end{cases}$$

- ▶ *hinge loss*: computationally tractable (convex) alternative

$$L(\hat{y}, y) = \max(0, 1 - y * \hat{y})$$

Hinge loss



Minimize hinge loss on training data, with regularization added

$$\operatorname{argmin}_{\mathbf{w}} \left\{ \underbrace{\sum_{i=1}^n \left(1 + \max(0, 1 - y * \mathbf{w}^T \mathbf{x}_i) \right)}_{\text{Training set error}} + \lambda \underbrace{\sum_{i=0}^d w_i^2}_{\text{Regularizer}} \right\}$$

Regularized logistic regression

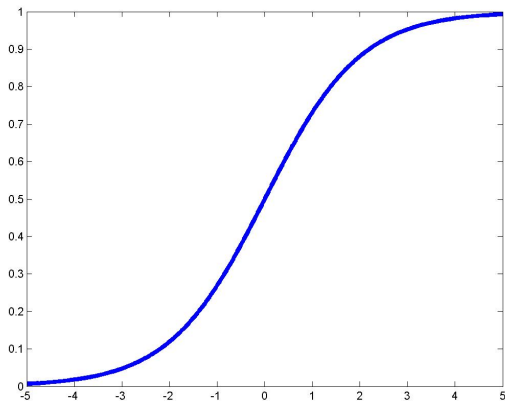
- ▶ **Given:** A set of data points $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ each of which belongs to one of the two classes denoted as -1 and 1 .
- ▶ **Desired:** A simple description of the function $p(y = 1|\mathbf{x})$, probability that instance belongs to class 1 given features \mathbf{x} (obviously, $p(y = -1|\mathbf{x}) = 1 - p(y = 1|\mathbf{x})$)
- ▶ **Approach:** Describe p by a logistic function:

$$p(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

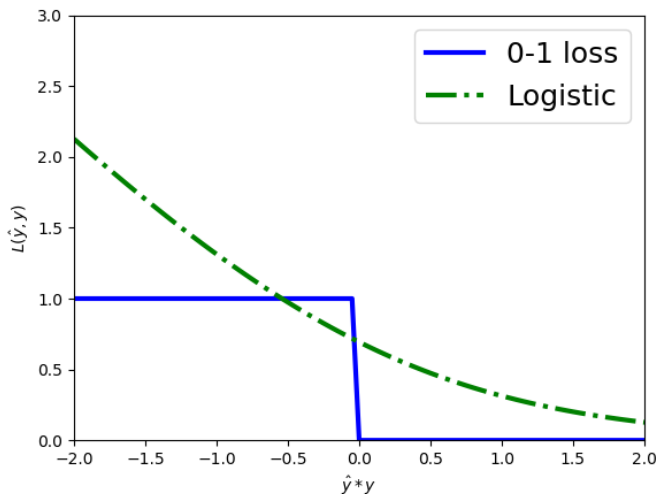
- ▶ A linear model, with logistic function used to squeeze the predictions between 0 and 1. How to learn the coefficients from data?

Sigmoid function

$$\frac{1}{1 + e^{-z}}$$



Logistic loss



Regularized logistic regression

Minimize the logarithm of the likelihood of training data, with regularization added

$$\operatorname{argmin}_{\mathbf{w}} \left\{ \underbrace{\sum_{i=1}^n \log \left(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i} \right)}_{\text{Training set error}} + \lambda \underbrace{\sum_{i=0}^d w_i^2}_{\text{Regularizer}} \right\}$$

Regularized logistic regression

- ▶ logistic regression: a very classical linear model for classification
- ▶ predictions scaled between 0 and 1, can be interpreted as probabilities
- ▶ regularization can be used to help control model complexity
- ▶ may lead to better classification accuracy than using standard ridge regression
- ▶ can be trained with standard gradient descent optimization
- ▶ `sklearn.linear_model.LogisticRegression`

Predicting multiple classes with binary classifiers

- ▶ Binary classifiers solve two-class classification problems
- ▶ How to solve problems with $k > 2$ classes?
- ▶ Multi-class classification: predict exactly one class
 - ▶ "This animal is a dog"
- ▶ Multi-label classification: predict a subset of classes
 - ▶ "This picture contains a dog and a cat"
- ▶ Option 1: some methods naturally suited to multiclass / multilabel problems
 - ▶ e.g. nearest neighbors, neural network with k outputs
- ▶ Option 2: reduction to binary classification
 - ▶ metaclassifier that can use any binary classification algorithm
 - ▶ e.g. logistic regression, support vector machine, ridge regression based classification

Reduction of multi-class to binary classification

| ID | Dog | Cat | Rabbit | Hamster |
|----------|-----|-----|--------|---------|
| Image #1 | 1 | 0 | 0 | 0 |
| Image #2 | 1 | 0 | 0 | 0 |
| Image #3 | 0 | 0 | 1 | 0 |
| Image #4 | 0 | 0 | 0 | 1 |
| Image #5 | 0 | 1 | 0 | 0 |

- ▶ Species recognition task: k different possible classes
- ▶ Images #1 and #2 contain dogs, #3 a rabbit etc.
- ▶ Multiclass data: each image belongs to exactly one class
- ▶ One-vs-all: use one-hot-encoding for the classes
- ▶ Train separate binary classifier for each column
- ▶ For prediction, choose class with highest predicted probability/confidence

Reduction of multi-label to binary classification

| ID | Dog | Cat | Rabbit | Hamster |
|----------|-----|-----|--------|---------|
| Image #1 | 1 | 1 | 0 | 0 |
| Image #2 | 1 | 0 | 0 | 1 |
| Image #3 | 0 | 0 | 0 | 0 |
| Image #4 | 0 | 0 | 0 | 1 |
| Image #5 | 1 | 1 | 1 | 0 |

- ▶ Similar, but now instances may belong to multiple classes simultaneously, or none at all
- ▶ Binary relevance method: train separate binary classifier for each label
- ▶ Prediction: choose e.g. all classes with > 0.5 predicted probability (more generally $> t$ predicted confidence score for some threshold t)

Non-linear models

- ▶ what if the model to be learned is highly non-linear, and nearest neighbour does not work well enough?
- ▶ kernel methods: generalize basic linear models to handling non-linear data
- ▶ neural networks models
- ▶ random forests, gradient boosting
- ▶ more complicated models, same basic idea: balance training error with model complexity
- ▶ outside the scope of this course

Ridge regression and kernel trick

