

## Lecture 1: Structural Risk Minimization

Instructor: Marion Neumann

Reading: FCML Ch1 (Linear Modeling); ESL 3.4.3, 10.6

## Learning Objective

Understand that many machine learning algorithms solve the **structural risk minimization problem**, which is essentially **minimizing a combination of loss function and model complexity penalty**.

## Application

Build a **spam filter** that works well on the *training data* and generalizes well to unseen *test data*.

In fact, this will be our first implementation project for the course. Take some time to answer the following warm-up questions:

(1) *How does our data look like?* (2) *What are the features?* (3) *What is the prediction task?* (4) *How well do you think a linear classifier will perform?* (5) *How do you measure the performance of a (linear) classifier?*



## 1 Introduction

## 1.1 Machine Learning Problem

Assume we have a **dataset**

$$D = \{(\mathbf{x}_i, y_i)\}_{i=1, \dots, n}, \quad (1)$$

our goal is to learn

$$\mathbf{y} = h(\mathbf{x}) \quad (2)$$

such that:

$$\begin{aligned} h(\mathbf{x}_i) &= y_i, \forall i = 1, \dots, n \text{ and} \\ h(\mathbf{x}^*) &= y^* \text{ for unseen test data} \end{aligned}$$

Note that we can write  $\mathbf{y} = h(f(\mathbf{x}))$  with  $f: \mathbf{x} \rightarrow \mathbb{R}$  and we get **classification models** with class labels  $\{-1, +1\}$  by choosing the sign function  $h(a) = \text{sign}(a)$  and **regression models** by using the identity function  $h(a) = I(a)$ .

Question: How do we find such a function  $h$ ?

Answer: Optimize some performance measure of the model (aka the function/hypothesis  $h$ ).

$\Rightarrow$  **minimize expected risk:**

$$\min_h R[h] = \min_h \int \underbrace{l(h(\mathbf{x}), \mathbf{y})}_{\text{e.g. squared loss}} dp(\mathbf{x}, \mathbf{y}), \quad (3)$$

where  $p(\mathbf{x}, \mathbf{y})$  is the joint probability of the data, which is **unknown**.

⇒ use empirical risk instead:

$$\min_h R_{emp}[h] = \min_h \frac{1}{n} \sum_{i=1}^n l(h(\mathbf{x}_i), y_i) \quad (4)$$

*Empirical risk minimization* minimizes the training error. However, this tends to *overfit* to the training data, and typically a simpler model is preferred (*Occam's razor*).

**Recap:** *Occam's razor* says one should pick the simplest model that adequately explains the data.

## 1.2 Structural Risk Minimization

The goal of *structural risk minimization* is to balance fitting the training data against model complexity. **Training** means then to learn the model parameters by solving the following optimization problem:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) = \min_{\mathbf{w}} \underbrace{\frac{1}{n} \sum_{i=1}^n l(h_{\mathbf{w}}(\mathbf{x}_i), y_i)}_{\text{training loss}} + \lambda \underbrace{r(\mathbf{w})}_{\text{regularizer}} \quad (5)$$

where the objective function  $\mathcal{L}(\mathbf{w})$  combines a *loss function* penalizing a high training error and a *regularizer* penalizing the model complexity.  $\lambda$  is a model *hyperparameter* that controls the trade-off between the terms. We are interested in choosing  $\lambda$  to minimize the *true risk* (test error). As the true risk is unknown we may resort to *cross-validation* to learn  $\lambda$ .<sup>1</sup> The science behind finding an ideal loss function is known as **Empirical Risk Minimization** (ERM). Extending the objective function to incorporate a regularizer leads to **Structural Risk Minimization** (SRM).

This provides us with a unified view on many machine learning methods. By plugging in different (surrogate) loss functions and regularizers, we obtain different machine learning models. Remember for example the **unconstrained SVM formulation**:

$$\min_{\mathbf{w}} C \underbrace{\sum_{i=1}^n \max[1 - y_i(\underbrace{\mathbf{w}^T \mathbf{x}_i + b}_{=f_{\mathbf{w}}(\mathbf{x}_i)}, 0]}_{\text{hinge-loss}} + \underbrace{\|\mathbf{w}\|_2^2}_{l_2\text{-regularizer}} \quad (6)$$

**SVM** uses the **hinge loss** as error measure and the  **$l_2$ -regularizer** to penalize complex solutions;  $\lambda = \frac{1}{C}$ .

<sup>1</sup>In Bayesian machine learning, we truly incorporate the choice of  $\lambda$  into the learning objective.

## 2 Loss Functions

### 2.1 Commonly Used Binary Classification Loss Functions

For binary classification we naturally want to minimize the zero-one loss  $l(h_{\mathbf{w}}(\mathbf{x}), y)$  also called **true classification error**. However, due to its **non-continuity** it is impractical to optimize and we resort to using various *surrogate loss functions*  $l(f_{\mathbf{w}}(\mathbf{x}), y)$ . Table 1 summarizes the most commonly used classification losses.

Table 1: loss functions for classification  $y \in \{-1, +1\}$

Loss $l(f_{\mathbf{w}}(\mathbf{x}_i), y_i)$	Usage	Comments
<b>Zero-one Loss:</b> $\delta(h_{\mathbf{w}}(\mathbf{x}_i) \neq y_i)$	true classification loss	Non-continuous and thus impractical to optimize.
<b>Hinge-Loss:</b> $\max[1 - f_{\mathbf{w}}(\mathbf{x}_i) * y_i, 0]^p$	<ul style="list-style-type: none"> <li>standard SVM (<math>p = 1</math>)</li> <li>(differentiable) squared hinge loss SVM (<math>p = 2</math>)</li> </ul>	When used for standard SVM, the loss function denotes margin length between linear separator and its closest point in either class. Only differentiable everywhere with $p = 2$ .
<b>Log-Loss:</b> $\log(1 + e^{-f_{\mathbf{w}}(\mathbf{x}_i)y_i})$	logistic regression	One of the most popular loss functions in machine learning, since its outputs are very well-tuned.
<b>Exponential Loss:</b> $e^{-f_{\mathbf{w}}(\mathbf{x}_i)y_i}$	AdaBoost	This function is very aggressive and thus sensitive to label noise. The loss of a mis-prediction increases exponentially with the value of $-f_{\mathbf{w}}(\mathbf{x}_i)y_i$ .

What do all these loss functions look like? The Illustration below shows the zero-one and exponential losses, where the input/x-axis is the "correctness" of the prediction  $f_{\mathbf{w}}(\mathbf{x}_i)y_i$ .

**Exercise 2.1.** Add the *hinge loss* and *log loss functions* to this plot.

#### Additional Notes on Classification Loss Functions

- Zero-one loss is zero when the prediction is correct, and one when incorrect.
- As  $z \rightarrow \infty$ , log-loss, exp-loss, and hinge loss become increasingly parallel.
- The exponential loss and the hinge loss are both upper bounds of the zero-one loss. (For the exponential loss, this is an important aspect in Adaboost.)

**Exercise 2.2.** In what scenario would you want to use the **Huber loss** instead of squared or **absolute loss**?

### 2.2 Commonly Used Regression Loss Functions

Unsurprisingly, regression models (where the predictions are reals) also have their own loss functions summarized in Table 2. Note that for regression the error is quantified as  $h_{\mathbf{w}}(\mathbf{x}) = f_{\mathbf{w}}(\mathbf{x})$ .

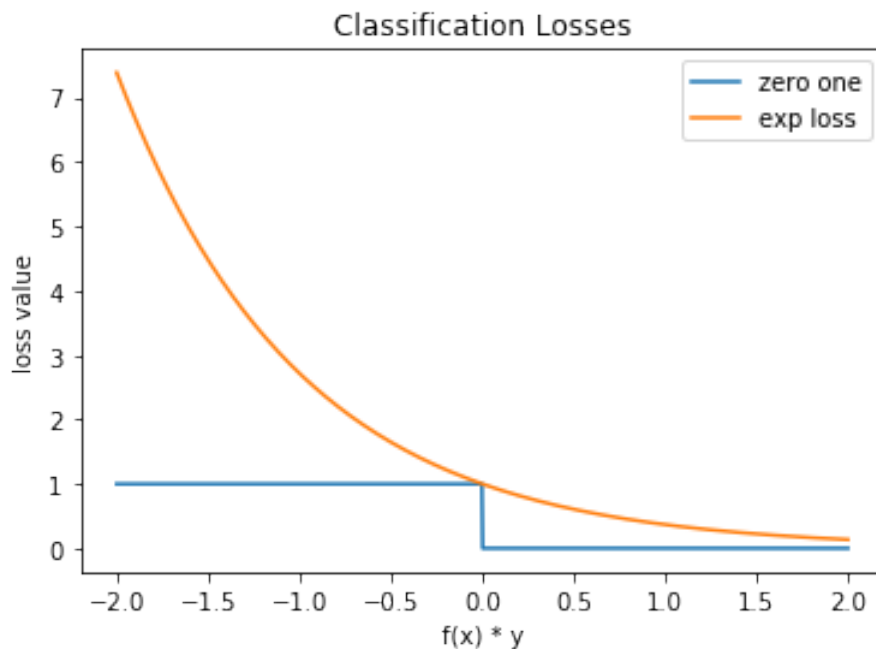


Table 2: Loss Functions With Regression,  $y \in \mathbb{R}$

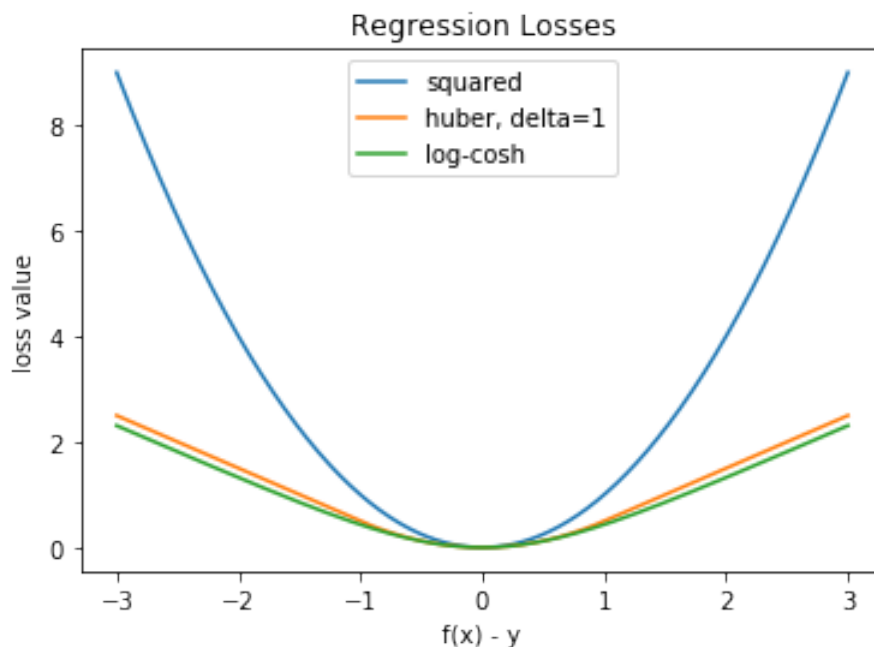
Loss $l(f_{\mathbf{w}}(\mathbf{x}_i), y_i)$	Comments
<p><b>Squared Loss:</b></p> $(f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2$	<ul style="list-style-type: none"> <li>• most popular regression loss function</li> <li>• <math>\mathbf{w}^*</math> will be related to the <b>mean</b> observations in <math>D^2</math></li> <li>• ADVANTAGE: differentiable everywhere</li> <li>• DISADVANTAGE: tries to accommodate every sample <math>\rightarrow</math> sensitive to outliers/noise</li> <li>• used in <b>Ordinary Least Squares (OLS)</b></li> </ul>
<p><b>Absolute Loss:</b></p> $ f_{\mathbf{w}}(\mathbf{x}_i) - y_i $	<ul style="list-style-type: none"> <li>• also a very popular loss function</li> <li>• <math>\mathbf{w}^*</math> will be related to the <b>median</b> observations in <math>D^3</math></li> <li>• ADVANTAGE: less sensitive to noise</li> <li>• DISADVANTAGE: <b>not differentiable at 0</b></li> </ul>

<sup>2</sup>For the same input location and multiple observations, the minimal squared loss is achieved for the mean of the observations.

<sup>3</sup>For the same input location and multiple observations, the minimal absolute loss is achieved for the median of the observations.

<p><b>Huber Loss:</b></p> $\begin{cases} \frac{1}{2} z_i^2 & \text{if }  z_i  < \delta \\ \delta( z_i  - \frac{\delta}{2}) & \text{otherwise} \end{cases}$ <p>where <math>z_i = f_{\mathbf{w}}(\mathbf{x}_i) - y_i</math></p>	<ul style="list-style-type: none"> <li>• also known as <b>Smooth Absolute Loss</b></li> <li>• once-differentiable</li> <li>• ADVANTAGE: “Best of Both Worlds” of <u>squared</u> and <u>absolute</u> loss</li> <li>• Takes on behavior of squared loss when loss is small, and absolute loss when loss is large.</li> </ul>
<p><b>Log-Cosh Loss:</b></p> $\log(\cosh(f_{\mathbf{w}}(\mathbf{x}_i) - y_i)),$ <p>where <math>\cosh(x) = \frac{e^x + e^{-x}}{2}</math></p>	<ul style="list-style-type: none"> <li>• ADVANTAGE: similar to Huber loss, but twice differentiable everywhere.</li> </ul>
<p><b><math>\varepsilon</math>-Insensitive Loss</b></p> $\begin{cases} 0 & \text{if }  z_i  < \varepsilon \\  z_i  - \varepsilon & \text{otherwise} \end{cases}$ <p>where <math>z_i = f_{\mathbf{w}}(\mathbf{x}_i) - y_i</math></p>	<ul style="list-style-type: none"> <li>• yields sparse solution (cf. <i>support vectors</i>)</li> <li>• used in SVM regression</li> <li>• <math>\varepsilon</math> regulates the sensitivity of the loss and hence the number of support vectors used in the SVM</li> </ul>

What do all these loss functions look like? The Illustration below shows the squared loss, huber loss with  $\delta = 1$ , and log-cosh loss, where the input/x-axis is the “correctness” of the prediction  $z = f_{\mathbf{w}}(\mathbf{x}_i) - y_i$ .



**Exercise 2.3.** Add the functions for the *absolute loss* and the *huber loss* with  $\delta = 2$  to this plot.

**Exercise 2.4.** In the same or a new plot, graph the functions for the  *$\varepsilon$ -insensitive loss* for  $\varepsilon = 1$  and  $\varepsilon = 0.1$ .

### 3 Regularizers

Similar to the SVM primal derivation, we can establish the following equivalency:

$$\min_{\mathbf{w}, b} \sum_{i=1}^n l(\mathbf{w}^T \mathbf{x}_i + b, y_i) + \lambda r(\mathbf{w}) \iff \min_{\mathbf{w}, b} \sum_{i=1}^n l(\mathbf{w}^T \mathbf{x}_i + b, y_i) \quad (7)$$

subject to:  $r(\mathbf{w}) \leq B$

For each  $\lambda \geq 0$ , there exists  $B \geq 0$  such that the two formulations in Eq. (7) are equivalent, and vice versa. In previous sections, the  $l_2$ -regularizer has been introduced as the component in SVM that reflects the complexity of solutions. Besides the  $l_2$ -regularizer, other types of useful regularizers and their properties are listed in Table 3.

Table 3: Types of regularizers

Regularizers $r(\mathbf{w})$	Properties
<b><math>l_2</math>-regularizer:</b> $r(\mathbf{w}) = \mathbf{w}^T \mathbf{w} = (\ \mathbf{w}\ _2)^2$	<ul style="list-style-type: none"> <li>• ADVANTAGE: strictly convex</li> <li>• ADVANTAGE: differentiable</li> <li>• DISADVANTAGE: uses weights on all features, i.e. relies on all features to some degree (ideally we would like to avoid this) - these are known as <u>dense solutions</u>.</li> </ul>
<b><math>l_1</math>-regularizer:</b> $r(\mathbf{w}) = \ \mathbf{w}\ _1$	<ul style="list-style-type: none"> <li>• convex (but not strictly)</li> <li>• DISADVANTAGE: not differentiable at 0 (the point which minimization is intended to bring us to)</li> <li>• ADVANTAGE: <u>sparse</u> (i.e. not <u>dense</u>) solutions</li> </ul>
<b>Elastic Net:</b> $\alpha \ \mathbf{w}\ _1 + (1 - \alpha) \ \mathbf{w}\ _2^2, \alpha \in [0, 1)$	<ul style="list-style-type: none"> <li>• ADVANTAGE: strictly convex (i.e. unique solution)</li> <li>• DISADVANTAGE: non-differentiable</li> </ul>
<b>lp-Norm, often <math>0 &lt; p \leq 1</math>:</b> $\ \mathbf{w}\ _p = \left( \sum_{i=1}^d  \mathbf{w} ^p \right)^{\frac{1}{p}}$	<ul style="list-style-type: none"> <li>• DISADVANTAGE: non-convex <math>\rightarrow</math> initialization dependent</li> <li>• ADVANTAGE: very sparse solutions</li> <li>• DISADVANTAGE: not differentiable</li> </ul>

Figure 1a shows plots of some common used regularizers. Note that regularizers are functions of  $\mathbf{w}$  and not  $\mathbf{x}_i$ . Figure 1b shows the effect that adding a regularizer to the loss function minimization has on the optimization problem and its solution.

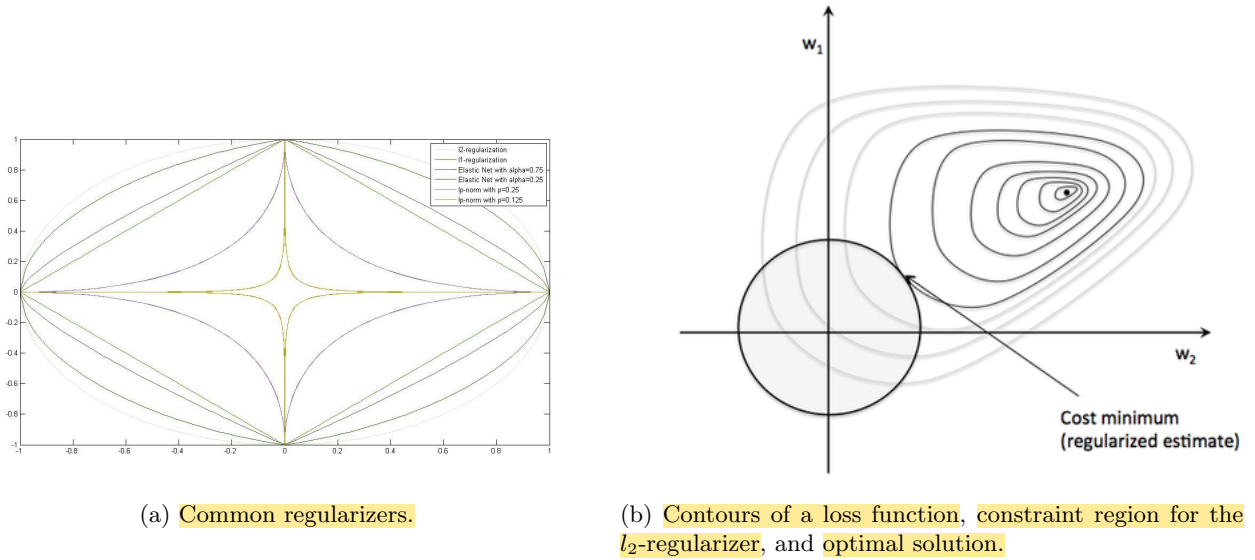


Figure 1: Illustrations for a two dimensional feature space ( $d = 2$ ).

### Exercise 3.1. Regularizers

- Add the constraint region for the  $l_1$  and elastic net regularizers to the plot in Figure 1b.
- The  $l_p$ -regularizer yields sparse solutions for  $\mathbf{w}$ . Explain why.
- What is the advantage of using the elastic net regularizer compared to  $l_1$  regularization? What is its advantage compared to  $l_2$  regularization? Briefly explain one advantage each.

## 4 Famous SRM models

This section includes several special cases of SRM – all preforming structural risk minimization, such as Ordinary Least Squares, Ridge Regression, Lasso, and Logistic Regression. Table 4 provides information on their loss functions, regularizers, as well as solutions.

**Notation:** Here we use  $X \in \mathbb{R}^{d \times n}$ , where  $d$  is the number of feature dimensions and  $n$  is the number of training data points<sup>4</sup>. **Caution:** this is the transpose of what they use in FCML.

Table 4: Special Cases of SRM

Loss and Regularizer	Properties	Solution
Ordinary Least Squares: $\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2$	<ul style="list-style-type: none"> <li>squared loss</li> <li>no regularization</li> </ul>	<ul style="list-style-type: none"> <li><math>\mathbf{w} = (X X^T)^{-1} X y^T</math></li> <li><math>X = [\mathbf{x}_1, \dots, \mathbf{x}_n]</math></li> <li><math>\mathbf{y} = [y_1, \dots, y_n]</math></li> </ul>

<sup>4</sup>I personally prefer this notation since you can quickly determine whether we are using inner or outer products. E.g. the inner product in matrix notation will then look like  $X^T X$ , which is more intuitive as it aligns with the inner product of vectors  $\mathbf{x}^T \mathbf{x}$ . However, both notations are found in the literature. So, always be careful about the definition of  $X$ .

<b>Ridge Regression:</b> $\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \ \mathbf{w}\ _2^2$	<ul style="list-style-type: none"> <li>• squared loss</li> <li>• <math>l_2</math>-regularization</li> </ul>	<ul style="list-style-type: none"> <li>• <math>\mathbf{w} = (XX^T + \lambda I)^{-1} X y^T</math></li> </ul>
<b>Lasso:</b> $\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \ \mathbf{w}\ _1$	<ul style="list-style-type: none"> <li>• + sparsity inducing (good for feature selection)</li> <li>• + convex</li> <li>• - not strictly convex (no unique solution)</li> <li>• - not differentiable (at 0)</li> </ul>	<ul style="list-style-type: none"> <li>• Solve with (sub)-gradient descent or LARS (least angle regression)</li> </ul>
<b>Logistic Regression:</b> $\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i(\mathbf{w}^T \mathbf{x}_i + b)})$	<ul style="list-style-type: none"> <li>• often also <math>l_1</math> or <math>l_2</math> regularized</li> </ul>	<ul style="list-style-type: none"> <li>• Solve with <b>gradient descent</b> or <b>Newton's method</b></li> <li>• <math>P(y = +1   \mathbf{x}) = \frac{1}{1 + e^{-y(\mathbf{w}^T \mathbf{x} + b)}}</math></li> </ul>
SVM: cf. Equation (6)	<ul style="list-style-type: none"> <li>• hinge-loss</li> <li>• <math>l_2</math>-regularization</li> </ul>	<ul style="list-style-type: none"> <li>• solve dual quadratic program (QP)</li> </ul>

## 5 Summary

SRM tries to find the best model by explicitly adding a weighted complexity penalty (regularization term) to the loss function optimization. **List of concepts and terms to understand from this lecture:**

- ERM
- SRM
- *overfitting*
- *underfitting*
- *Occam's razor*
- *training loss, testing loss*
- *surrogate loss*
- *regularizer*
- *hyperparameter*



### Exercise 5.1. Practice Retrieving!

For this summary exercise, it is intended that your answers are based on **your own** (current) understanding of the concepts (and not on the definitions you read and copy from these notes or from elsewhere). Don't hesitate to **say it out loud** to your seat neighbor, your pet or stuffed animal, or to yourself before **writing it down**. Research studies show that this practice of retrieval and phrasing out loud will help you retain the knowledge!

- (a) Using *your own words*, summarize each of the above concepts in 1-2 sentences by retrieving the knowledge from the top of your head.
- (b) What is the difference between ERM and SRM?
- (c) What is the difference between *overfitting* and *underfitting*?

And always remember: It's not bad to get it wrong. *Getting it wrong is part of learning!* Use your notes or other resources to get the correct answer or come to our office hours to get help!

## Our Application



With respect to our **spam filter application**, we are now able to come up with objective functions, aka learning models, that will (hopefully) work well on the training data and are able to generalize well to unseen test data.

The next step is to actually build the spam filter, which means we will need to solve the SRM problem (Eq. 5) for various combinations of loss functions and regularizers. In the next lecture, we will cover a variety of optimization techniques to help us with that.