

Introduction to Machine Learning

Lecture II

Oleg Filatov ¹ Stepan Zakharov ²

¹DESY ²Novosibirsk State University



Outline

- **Linear models: regression**
 - Analytic solution
 - (Stochastic) Gradient Descent
 - When to stop training?
 - (Under) Overtraining and its treatment
 - Loss functions & metrics
- **Linear models: classification**
 - Binary classifier
 - Logistic regression
 - Metrics
 - Multiclass case

Recap — SL fundamentals

Given:

- Set of objects that forms a training sample $\{x_1 \dots x_N\} \subset X$
- Set of answers $y_i = g(x_i)$, $\{y_1 \dots y_N\} \subset Y$

Find:

- $a: X \rightarrow Y$ – algorithm a that approximates g on the whole set X

Define:

- Loss function $\mathcal{L}(a(x), y)$
- $Q(a, X^N) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(a(x_i), y_i)$

$$\text{Then: } \hat{a} = \arg \min_{a \in A} Q(a, X^N)$$

Recap — ○ features and answers

Objects are described by a set of
features $f_j(x) \in D_j$:

- $D_j = \{0, 1\} \Rightarrow$ binary feature
- $|D_j| < \infty \Rightarrow$ categorical feature
- $|D_j| < \infty + \text{total order} \Rightarrow$ ordinal feature
- $D_j = \mathbb{R} \Rightarrow$ numerical feature

Type of answer $y_i \in Y$ generally defines
problem type:

- $Y = \{0, 1\} \Rightarrow$ binary classification
- $|Y| < \infty \Rightarrow$ multi-class classification
- $|Y| < \infty + \text{total order} \Rightarrow$ ranking
- $Y = \mathbb{R} \Rightarrow$ regression

Linear models

regression

Linear regression

- $Y = \mathbb{R}$
- N objects with K real features: $D = \mathbb{R}^K$
- $a(\mathbf{x}, \boldsymbol{\theta}) = \theta_0 + \sum_{j=1}^K f_j(\mathbf{x}) \cdot \theta_j$
- Extend and reassign:
 $[1, f_1(\mathbf{x}), \dots, f_K(\mathbf{x})] \equiv \mathbf{x}$
 $[\theta_0, \dots, \theta_K] \equiv \boldsymbol{\theta}$
- Then $a(\mathbf{x}, \boldsymbol{\theta}) = \langle \mathbf{x}, \boldsymbol{\theta} \rangle$
- Minimization problem:
 - $\mathcal{L}(\boldsymbol{\theta}) = (\langle \mathbf{x}, \boldsymbol{\theta} \rangle - y)^2$
 - $Q(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N (\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle - y_i)^2$
- Then we need to minimize $Q(\boldsymbol{\theta})$ by varying $\boldsymbol{\theta}$:
 - $\hat{\mathbf{a}} = \arg \min_{\boldsymbol{\theta} \in \Theta} Q(\boldsymbol{\theta})$

Solution can be found analytically

After defining feature matrix as:

$$F = \begin{bmatrix} f_1(\mathbf{x}_1) & \cdots & f_K(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ f_1(\mathbf{x}_N) & \cdots & f_K(\mathbf{x}_N) \end{bmatrix}$$

Minimization problem for linear regression:

- $Q(\boldsymbol{\theta}) = \frac{1}{N}(F\boldsymbol{\theta} - \mathbf{y})^T(F\boldsymbol{\theta} - \mathbf{y}) = \boldsymbol{\theta}^T F^T F \boldsymbol{\theta} - 2\mathbf{y}^T F \boldsymbol{\theta} + \mathbf{y}^T \mathbf{y}$
- $\frac{\partial Q}{\partial \boldsymbol{\theta}^T} = 2F^T F \boldsymbol{\theta} - 2F^T \mathbf{y} = 0$

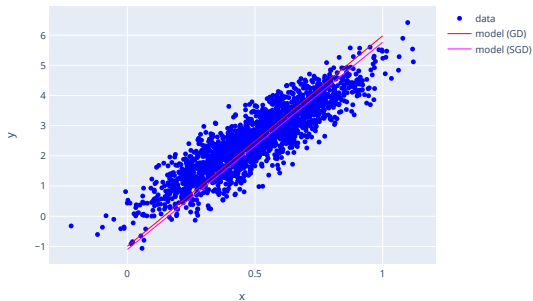
$$\hat{\boldsymbol{\theta}} = (F^T F)^{-1} F^T \mathbf{y}$$

Problem 1: calculation of $(F^T F)^{-1}$ has **complexity of $O(K^5 N)$** \Rightarrow computationally-intensive for large number of features

Problem 2: F is often **ill-conditioned** \Rightarrow computational errors causes instabilities and extremely large solutions

Better approach? —○ example

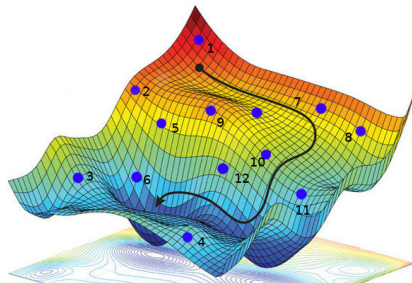
- Let's consider a problem where objects have only one feature:
 $\mathbf{x} = [1, f_1(\mathbf{x})]$
- Then $\boldsymbol{\theta} = [\theta_0, \theta_1]$
- How can we find minimum of $Q(\boldsymbol{\theta})$ without solving matrix equation?



Better approach? —○ random search

- Start with some vector θ^0
- Initialize $\hat{\theta} = \theta^0$
- ↓ Take randomly another vector θ'
- ↓ Compare $Q(\theta')$ and $Q(\hat{\theta})$
- ↓ $\hat{\theta} = \arg \min_{\hat{\theta}, \theta'} (Q(\hat{\theta}), Q(\theta'))$
- ↓ Repeat until convergence

Problem: random search of $\hat{\theta}$ is not effective for practical use in high dimensional Θ space

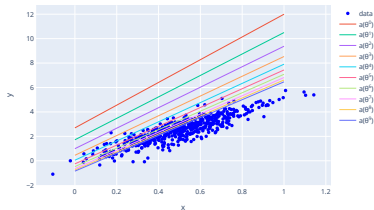
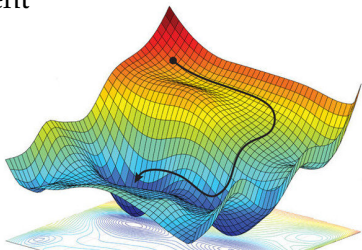


Q as a function of θ_0 and θ_1 .
Blue dots represent Q values for different θ

Better approach? —○ gradient descent

Idea: bring up additional information about the direction to the minimum of Q using its **antigradient**

- Start with some vector θ^0
- $\hat{\theta} \leftarrow \theta^0$
- ↓ Calculate direction to the minimum:
$$-\nabla Q(\hat{\theta}) = -\left[\frac{\partial Q}{\partial \theta_0}, \frac{\partial Q}{\partial \theta_1}\right] \hat{\theta}$$
- ↓ Move towards the minimum:
$$\hat{\theta} \leftarrow \hat{\theta} - \eta \nabla Q(\hat{\theta})$$
- ↓ Repeat until convergence

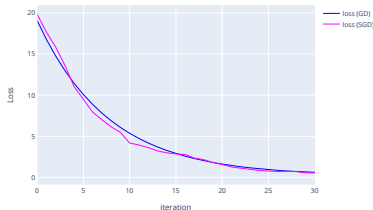


Better approach? —○ stochastic gradient descent

Gradient descent

Calculating $-\nabla Q$ over all objects

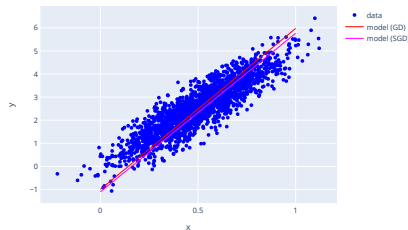
- ✓ Precise gradient
- ✗ Can be time-consuming for large dataset



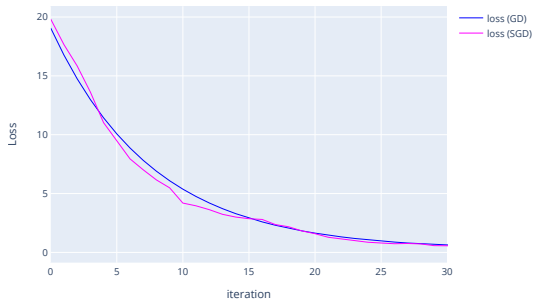
Stochastic gradient descent

Calculating $-\nabla Q$ over minibatch

- ✓ Big data friendly
- ✗ Noisy gradient



When to stop training?



When to stop training? — vanilla approaches

1 Stop once reached predefined number of iteration

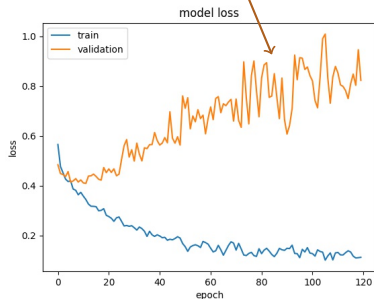
? what if too early? → **undertraining**

? what if too late? → **overtraining**

2 Stop once training loss is not improving

- can use $\mathcal{L}(a(x), y)$ as an error estimator
 - measure it on training data at each iteration
 - stop once $\mathcal{L}_{\text{train}}$ is not improving by predefined ε
- monitor the progress online throughout the training
- can overfit easily

badly overfitted but $\mathcal{L}_{\text{train}}$ is decreasing

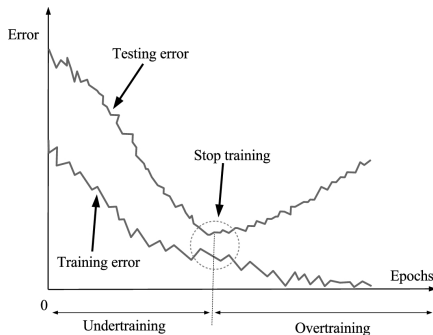


Note: you can use **any metric** instead of \mathcal{L}

When to stop training? —○ using a test set

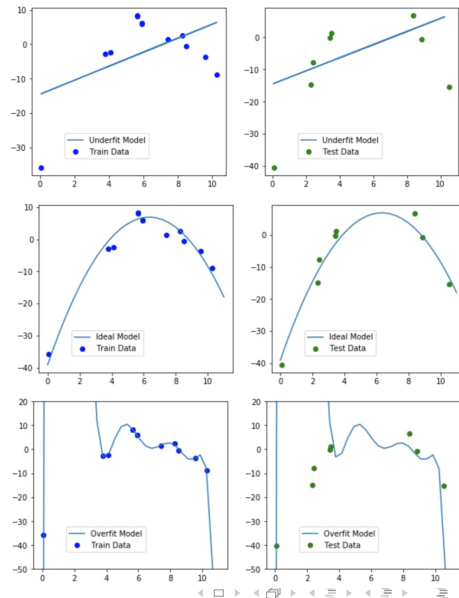
- To check model's capability to generalise we apply it to **unseen test data**
 - Can do this not only after training is done but **throughout**
 - So check train and test error for each iteration
 - And stop once test metric/loss is not improving by $\varepsilon \rightarrow$ **early stopping**
- detect and prevent overfitting before it happens

Note: "test set" and "validation set" are often used interchangeably



Example — ○ measuring model performance

- After training is done, we want to measure model's performance by applying it to test (unseen) data
- In this example, upper model performs bad both on train and test sets
- Center model shows acceptable results on train and test sets
- Lower model performs well on train and fails on test

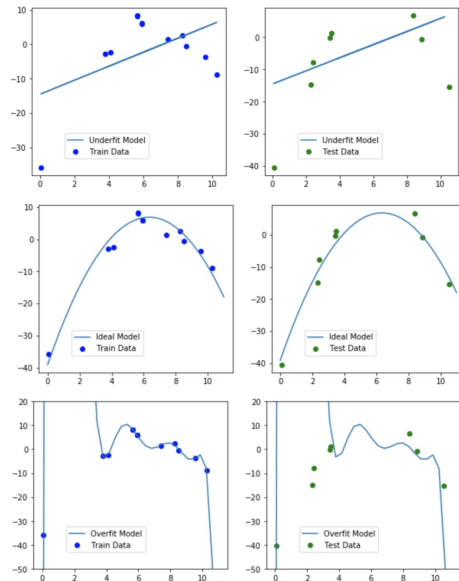


Example — measuring model performance

- After training is done, we want to measure model's performance by applying it to test (unseen) data
- In this example, upper model performs bad both on train and test sets
- Center model shows acceptable results on train and test sets
- Lower model performs well on train and fails on test



Let's analyze these outcomes



Possible outcomes

Train – bad
Test – bad



Undertraining

- Model is too simple
- Model can't capture the underlying structure of data
- Low variance, high bias

Train – good
Test – bad



Overtraining

- Model is too complicated
- Model captures noise instead of structure of data
- High variance, low bias

What to do?

- Increase model's capacity
- Train longer
- Catch it before it happens (e.g. early stopping)
- Reflect about noise level vs model complexity

Possible outcomes

Test – bad



Undertraining

- Model is too simple
- Model can't capture the underlying structure of data
- Low variance, high bias

Test – bad



Overtraining

- Model is too complicated
- Model captures noise instead of structure of data
- High variance, low bias

What to do?

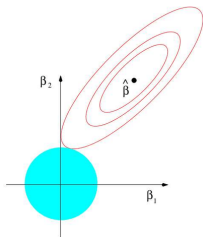
- Increase model's capacity
- Train longer
- Catch it before it happens (e.g. early stopping)
- Reflect about noise level vs model complexity
- **Automatically** limit algorithm space → **regularization**

Regularization

playgrounds [click](#), [click](#)

L2 regularization (Tikhonov)

$$Q(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N (\langle \mathbf{x}, \boldsymbol{\theta} \rangle - y_i)^2 + \lambda \sum_{j=1}^K \boldsymbol{\theta}_j^2$$

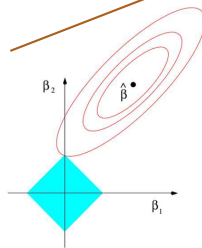


penalize model for too large weights

L1 regularization (LASSO)

least absolute shrinkage and selection operator

$$Q(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N (\langle \mathbf{x}, \boldsymbol{\theta} \rangle - y_i)^2 + \lambda \sum_{j=1}^K |\boldsymbol{\theta}_j|$$

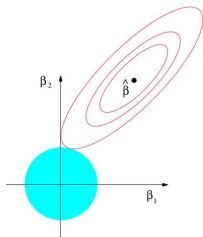


Regularization

playgrounds [click](#), [click](#)

L2 regularization (Tikhonov)

$$Q(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N (\langle \mathbf{x}, \boldsymbol{\theta} \rangle - y_i)^2 + \lambda \sum_{j=1}^K \boldsymbol{\theta}_j^2$$

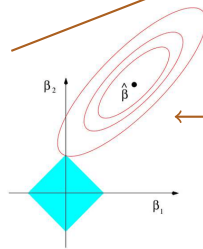


penalize model for too large weights

L1 regularization (LASSO)

least absolute shrinkage and selection operator

$$Q(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N (\langle \mathbf{x}, \boldsymbol{\theta} \rangle - y_i)^2 + \lambda \sum_{j=1}^K |\boldsymbol{\theta}_j|$$



autozeroing
uninformative
weights!

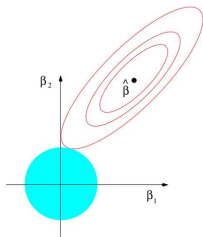
feature selection

Regularization

playgrounds [click](#), [click](#)

L2 regularization (Tikhonov)

$$Q(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N (\langle \mathbf{x}, \boldsymbol{\theta} \rangle - y_i)^2 + \lambda \sum_{j=1}^K \theta_j^2$$

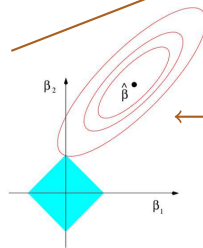


penalize model for too large weights

L1 regularization (LASSO)

least absolute shrinkage and selection operator

$$Q(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N (\langle \mathbf{x}, \boldsymbol{\theta} \rangle - y_i)^2 + \lambda \sum_{j=1}^K |\theta_j|$$



autozeroing
uninformative
weights!

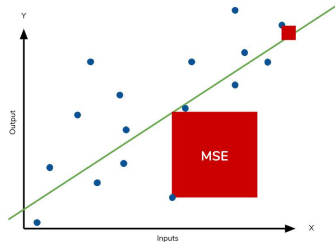
feature selection

λ is a **hyperparameter** \rightarrow more in homework

Why did we choose $\mathcal{L}(a(\mathbf{x}), y) = (a(\mathbf{x}) - y)^2$?

Why did we choose $\mathcal{L}(a(\mathbf{x}), y) = (a(\mathbf{x}) - y)^2$?

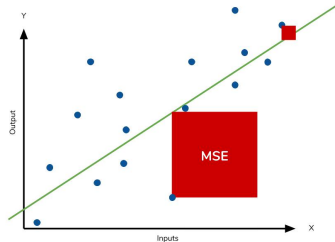
- we chose it because it represents model's error
- which we want to be as small as possible (keeping overfitting in mind)
- and therefore we **directly** minimize its **average over the dataset** in training



$$\mathcal{L}_{\text{train}} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(a(\mathbf{x}_i), y_i)$$

Why did we choose $\mathcal{L}(a(\mathbf{x}), y) = (a(\mathbf{x}) - y)^2$?

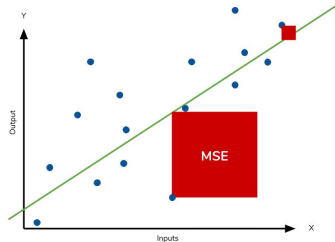
- we chose it because it represents model's error
- which we want to be as small as possible (keeping overfitting in mind)
- and therefore we **directly** minimize its **average over the dataset** in training



$$\mathcal{L}_{\text{train}} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(a(\mathbf{x}_i), y_i) \equiv \mathbb{E}_{p(\mathbf{x}, y)}[\mathcal{L}(a(\mathbf{x}), y)] \leftarrow \text{expected value}$$

Why did we choose $\mathcal{L}(a(x), y) = (a(x) - y)^2$?

- we chose it because it represents model's error
- which we want to be as small as possible (keeping overfitting in mind)
- and therefore we **directly** minimize its **average over the dataset** in training



$$\mathcal{L}_{\text{train}} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(a(x_i), y_i) \equiv \mathbb{E}_{p(x,y)}[\mathcal{L}(a(x), y)] \leftarrow \text{expected value}$$

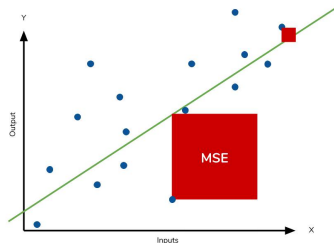
There is something more behind this...

Art of choosing loss function

—○ Mean Squared Error (L2)

$$\text{MSE}(a, X, Y) = \mathbb{E}_{p(x,y)} [(a(x) - y)^2]$$

- Grows quadratically with the error \Rightarrow punish for large mistakes \Rightarrow **sensitive to outliers** ("mean-unbiased")
- **Symmetrical** for under/overestimates
- **Smooth** \Rightarrow no problem with finding optimum
- **Doesn't preserve units** of measurement \Rightarrow poorly interpretable
- **Good statistical properties** \Rightarrow might give you BLUE (see [Gauss-Markov theorem](#))



Art of choosing loss function

—○ MSE variations

- $\text{RMSE} = \sqrt{\text{MSE}}$

→ preserve units of measurement \Rightarrow interpretable

- $\text{MSLE} = \mathbb{E}_{p(x,y)} [(\log(a(x) + 1) - \log(y + 1))^2]$

→ penalize on the "order of magnitude" scale

→ not symmetric for over/underestimation

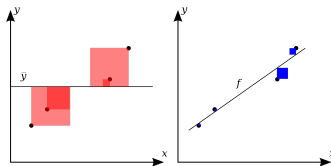
- **Coeff. of determination**

$$R^2 = 1 - \frac{\mathbb{E}_{p(x,y)} [(a(x) - y)^2]}{\mathbb{E}_{p(x,y)} [(\bar{y} - y)^2]}$$

→ fraction of "explained" variance: $1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$

→ "normalised" MSE

$$\text{MSE} = \mathbb{E}_{p(x,y)} [(a(x) - y)^2]$$



Art of choosing loss function

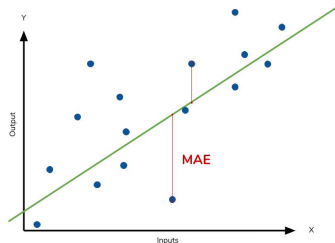
—○ Mean Absolute Error (L1)

$$\text{MAE}(a, X, Y) = \mathbb{E}_{p(x,y)} [|a(x) - y|]$$

- But for example, we could also take absolute value

* btw why not $\mathcal{L} = \mathbb{E}[a(x) - y]$?

- **Simple to interpret**
- **Symmetrical** for under/overestimates
- **Robust** to outliers ("median-unbiased")
- **Not smooth**
- **Not backed** much with statistics



Art of choosing loss function

—○ MAE variations

- $\text{MAPE} = \mathbb{E}_{\mathbf{p}(\mathbf{x},y)} \left[\left| \frac{y - a(\mathbf{x})}{y} \right| \right]$

→ "normalised" MAE

- $\text{SMAPE} = 100\% \cdot \mathbb{E}_{\mathbf{p}(\mathbf{x},y)} \left[\frac{|y - a(\mathbf{x})|}{|y| + |a(\mathbf{x})|} \right]$

→ symmetrized MAPE

→ not really symmetrized for over/underestimation

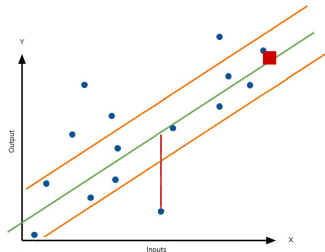
$$\text{MAE} = \mathbb{E}_{\mathbf{p}(\mathbf{x},y)} [|a(\mathbf{x}) - y|]$$

Art of choosing loss function

—○ Huber Loss

$$\text{HL}(a, X, Y) = \mathbb{E}_{P(x,y)} \left\{ \begin{array}{ll} \frac{1}{2} (a(x) - y)^2, & |a(x) - y| < \delta \\ \delta |a(x) - y| - \frac{1}{2} \delta^2, & \text{otherwise} \end{array} \right\}$$

- Or we could also combine them altogether
- **Robust** to outliers in tails
- **Better performance** near minimum (smaller variance)
- δ is a **hyperparameter** \Rightarrow should be tuned

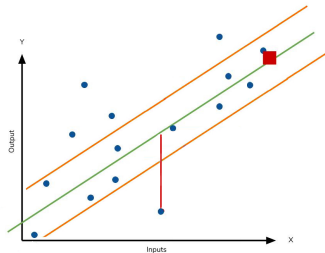


Art of choosing loss function

—○ Huber Loss

$$\text{HL}(a, X, Y) = \mathbb{E}_{P(x,y)} \left\{ \begin{array}{ll} \frac{1}{2} (a(x) - y)^2, & |a(x) - y| < \delta \\ \delta |a(x) - y| - \frac{1}{2} \delta^2, & \text{otherwise} \end{array} \right\}$$

- Or we could also combine them altogether
- **Robust** to outliers in tails
- **Better performance** near minimum (smaller variance)
- δ is a **hyperparameter** \Rightarrow should be tuned



All of the above losses can (and should) be used as metrics!

Linear models

classification

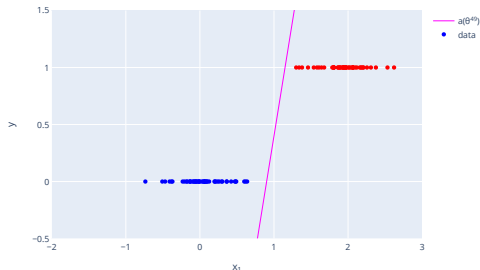
Classification

- $Y = \{-1, +1\}$
- N objects with K real features: $D = \mathbb{R}^K$
- $a(\mathbf{x}, \boldsymbol{\theta}) = \theta_0 + \sum_{j=1}^K f_j(\mathbf{x}) \cdot \theta_j$
- Extend and reassign:
 $[1, f_1(\mathbf{x}), \dots, f_K(\mathbf{x})] \equiv [1, x_1, \dots, x_K] \equiv \mathbf{x}$
 $[\theta_0, \dots, \theta_K] \equiv \boldsymbol{\theta}$
- Then $a(\mathbf{x}, \boldsymbol{\theta}) = \langle \mathbf{x}, \boldsymbol{\theta} \rangle$
- Minimization problem:
 - $\mathcal{L}(\boldsymbol{\theta}) = [\text{sign}\langle \mathbf{x}, \boldsymbol{\theta} \rangle \neq y]$
 - $Q(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N [\text{sign}\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle \neq y_i]$
- Then we need to minimize $Q(\boldsymbol{\theta})$ by varying $\boldsymbol{\theta}$:
 - $\hat{a} = \arg \min_{\boldsymbol{\theta} \in \Theta} Q(\boldsymbol{\theta})$

Example

- Let's consider the problem of two class (red, blue) classification
- We need to find a line* that provides the best separation of objects from different classes
- Objects have single feature:
 $\mathbf{x} = [1, x_1]$
- Then $\boldsymbol{\theta} = [\theta_0, \theta_1]$

*for a model with multiple features
we are looking for hyperplanes



Minimization problem for classification

- $\mathcal{L}(\theta) = [a(\mathbf{x}, \theta) \neq y] = [\text{sign}\langle \mathbf{x}, \theta \rangle \neq y] \leftarrow$ step function
- $Q(\theta) = \mathbb{E}_{p(\mathbf{x}, y)}[\text{sign}\langle \mathbf{x}, \theta \rangle \neq y] \leftarrow$ accuracy
- Then we need to minimize $Q(a(\theta))$ by varying θ

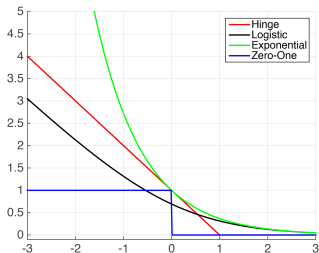
Problem 1: $\mathcal{L}(\theta)$ isn't a smooth function \Rightarrow gradient methods are not applicable

Problem 2: Discrete $\{-1, +1\}$ output from the classifier

From binary classification to logistic regression

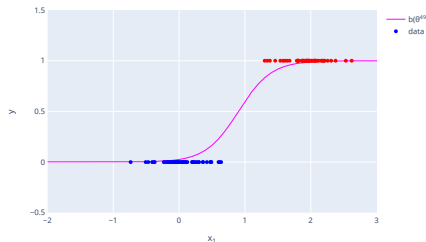
Problem 1:

Solution: use smooth approximations



Problem 2:

Solution: reframe the problem to probabilistic perspective
 \Rightarrow logistic regression



Understanding the concept of probability estimation

- Assume that algorithm $a(x, \theta)$ returns probability that object x_i is from positive class $y_i = +1$

Understanding the concept of probability estimation

- Assume that algorithm $a(x, \theta)$ returns probability that object x_i is from positive class $y_i = +1$
- Probability that x_i belongs to class y_i :
$$p_{\theta}(y_i|x_i) = a(x_i, \theta)^{[y_i=+1]} (1 - a(x_i, \theta))^{[y_i=-1]}$$

Understanding the concept of probability estimation

- Assume that algorithm $a(x, \theta)$ returns probability that object x_i is from positive class $y_i = +1$
- Probability that x_i belongs to class y_i :
$$p_{\theta}(y_i|x_i) = a(x_i, \theta)^{[y_i=+1]} (1 - a(x_i, \theta))^{[y_i=-1]}$$
- Function to be optimised: $\mathcal{L}(\theta) = \prod_{i=1}^N p_{\theta}(y_i|x_i) \leftarrow \text{maximize probability to observe such data}$

Understanding the concept of probability estimation

- Assume that algorithm $a(x, \theta)$ returns probability that object x_i is from positive class $y_i = +1$

- Probability that x_i belongs to class y_i :

$$p_{\theta}(y_i|x_i) = a(x_i, \theta)^{[y_i=+1]} (1 - a(x_i, \theta))^{[y_i=-1]}$$

- Function to be optimised: $\mathcal{L}(\theta) = \prod_{i=1}^N p_{\theta}(y_i|x_i)$ ← maximize probability to observe such data

- Product is not useful, let's make sum:

$$\log \mathcal{L}(\theta) = \sum_{i=1}^N [y_i = 1] \log a(x_i, \theta) + [y_i = -1] \log (1 - a(x_i, \theta))$$

Understanding the concept of probability estimation

- Assume that algorithm $a(x, \theta)$ returns probability that object x_i is from positive class $y_i = +1$
- Probability that x_i belongs to class y_i :
$$p_{\theta}(y_i|x_i) = a(x_i, \theta)^{[y_i=+1]} (1 - a(x_i, \theta))^{[y_i=-1]}$$
- Function to be optimised: $\mathcal{L}(\theta) = \prod_{i=1}^N p_{\theta}(y_i|x_i)$ ← maximize probability to observe such data
- Product is not useful, let's make sum:
$$\log \mathcal{L}(\theta) = \sum_{i=1}^N [y_i = 1] \log a(x_i, \theta) + [y_i = -1] \log (1 - a(x_i, \theta))$$
- Substitute targets: $-1 \rightarrow 0$

Understanding the concept of probability estimation

- Assume that algorithm $a(x, \theta)$ returns probability that object x_i is from positive class $y_i = +1$
- Probability that x_i belongs to class y_i :
$$p_{\theta}(y_i|x_i) = a(x_i, \theta)^{[y_i=+1]} (1 - a(x_i, \theta))^{[y_i=-1]}$$
- Function to be optimised: $\mathcal{L}(\theta) = \prod_{i=1}^N p_{\theta}(y_i|x_i)$ ← maximize probability to observe such data
- Product is not useful, let's make sum:
$$\log \mathcal{L}(\theta) = \sum_{i=1}^N [y_i = 1] \log a(x_i, \theta) + [y_i = -1] \log (1 - a(x_i, \theta))$$
- Substitute targets: $-1 \rightarrow 0$
- So that: $\log \mathcal{L}(\theta) = \sum_{i=1}^N y_i \log a(x_i, \theta) + (1 - y_i) \log (1 - a(x_i, \theta)) = \mathbb{E}_{p(x,y)} [\log(p_{\theta}(y|x))]$

Understanding the concept of probability estimation

- Assume that algorithm $a(x, \theta)$ returns probability that object x_i is from positive class $y_i = +1$
- Probability that x_i belongs to class y_i :
$$p_{\theta}(y_i|x_i) = a(x_i, \theta)^{[y_i=+1]} (1 - a(x_i, \theta))^{[y_i=-1]}$$
- Function to be optimised: $\mathcal{L}(\theta) = \prod_{i=1}^N p_{\theta}(y_i|x_i)$ ← maximize probability to observe such data
- Product is not useful, let's make sum:
$$\log \mathcal{L}(\theta) = \sum_{i=1}^N [y_i = 1] \log a(x_i, \theta) + [y_i = -1] \log (1 - a(x_i, \theta))$$
- Substitute targets: $-1 \rightarrow 0$ aka log-likelihood
↓
- So that: $\log \mathcal{L}(\theta) = \sum_{i=1}^N y_i \log a(x_i, \theta) + (1 - y_i) \log (1 - a(x_i, \theta)) = \mathbb{E}_{p(x,y)} [\log(p_{\theta}(y|x))]$

Understanding the concept of probability estimation

- Assume that algorithm $a(x, \theta)$ returns probability that object x_i is from positive class $y_i = +1$
- Probability that x_i belongs to class y_i :
$$p_{\theta}(y_i|x_i) = a(x_i, \theta)^{[y_i=+1]} (1 - a(x_i, \theta))^{[y_i=-1]}$$
- Function to be optimised: $\mathcal{L}(\theta) = \prod_{i=1}^N p_{\theta}(y_i|x_i)$ ← maximize probability to observe such data
- Product is not useful, let's make sum:
$$\log \mathcal{L}(\theta) = \sum_{i=1}^N [y_i = 1] \log a(x_i, \theta) + [y_i = -1] \log (1 - a(x_i, \theta))$$
- Substitute targets: $-1 \rightarrow 0$ aka log-likelihood
↓
- So that: $\log \mathcal{L}(\theta) = \sum_{i=1}^N y_i \log a(x_i, \theta) + (1 - y_i) \log (1 - a(x_i, \theta)) = \mathbb{E}_{p(x,y)} [\log(p_{\theta}(y|x))]$
- Finally: $\hat{a} = \arg \max_{\theta \in \Theta} \log \mathcal{L}(\theta)$

How to make classifier return probabilities?

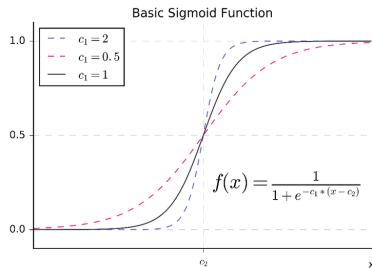
- $a(\mathbf{x}, \boldsymbol{\theta})$ has a range of \mathbb{R} , but it should be $[0, 1]$

How to make classifier return probabilities?

- $a(\mathbf{x}, \boldsymbol{\theta})$ has a range of \mathbb{R} , but it should be $[0, 1]$
- Let's think of a function that converts $\mathbb{R} \mapsto [0, 1]$

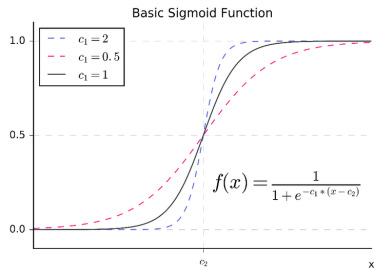
How to make classifier return probabilities?

- $a(\mathbf{x}, \boldsymbol{\theta})$ has a range of \mathbb{R} , but it should be $[0, 1]$
- Let's think of a function that converts $\mathbb{R} \mapsto [0, 1]$
- **Sigmoid function** looks nice: $\sigma(z) = \frac{1}{1 + e^{-z}}$



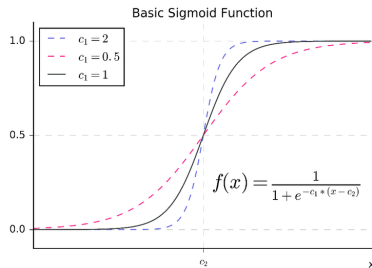
How to make classifier return probabilities?

- $a(\mathbf{x}, \boldsymbol{\theta})$ has a range of \mathbb{R} , but it should be $[0, 1]$
- Let's think of a function that converts $\mathbb{R} \mapsto [0, 1]$
- **Sigmoid function** looks nice: $\sigma(z) = \frac{1}{1 + e^{-z}}$
- Wrap it around: $a_{\sigma}(\mathbf{x}, \boldsymbol{\theta}) = \sigma(a(\mathbf{x}, \boldsymbol{\theta})) = \frac{1}{1 + e^{-\langle \mathbf{x}, \boldsymbol{\theta} \rangle}}$



How to make classifier return probabilities?

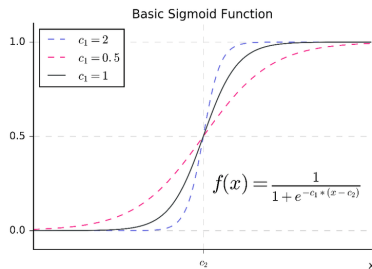
- $a(\mathbf{x}, \boldsymbol{\theta})$ has a range of \mathbb{R} , but it should be $[0, 1]$
- Let's think of a function that converts $\mathbb{R} \mapsto [0, 1]$
- **Sigmoid function** looks nice: $\sigma(z) = \frac{1}{1 + e^{-z}}$
- Wrap it around: $a_{\sigma}(\mathbf{x}, \boldsymbol{\theta}) = \sigma(a(\mathbf{x}, \boldsymbol{\theta})) = \frac{1}{1 + e^{-\langle \mathbf{x}, \boldsymbol{\theta} \rangle}}$
- Now it is possible to plug a_{σ} into log-likelihood and train:



How to make classifier return probabilities?

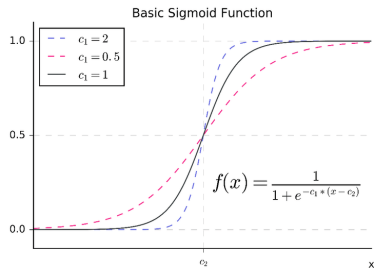
- $a(\mathbf{x}, \boldsymbol{\theta})$ has a range of \mathbb{R} , but it should be $[0, 1]$
- Let's think of a function that converts $\mathbb{R} \mapsto [0, 1]$
- **Sigmoid function** looks nice: $\sigma(z) = \frac{1}{1 + e^{-z}}$
- Wrap it around: $a_\sigma(\mathbf{x}, \boldsymbol{\theta}) = \sigma(a(\mathbf{x}, \boldsymbol{\theta})) = \frac{1}{1 + e^{-\langle \mathbf{x}, \boldsymbol{\theta} \rangle}}$
- Now it is possible to plug a_σ into log-likelihood and train:

$$\begin{aligned}\hat{a} &= \arg \max_{\boldsymbol{\theta} \in \Theta} \log \mathcal{L}(\boldsymbol{\theta}) = \\ &= \arg \max_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^N y_i \log \sigma(\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle) + (1 - y_i) \log (1 - \sigma(\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle))\end{aligned}$$



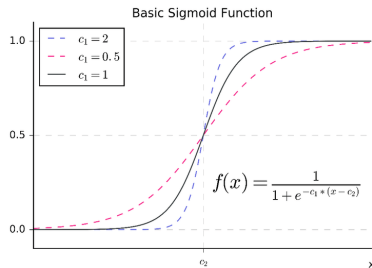
How to make classifier return probabilities?

- $a(\mathbf{x}, \boldsymbol{\theta})$ has a range of \mathbb{R} , but it should be $[0, 1]$
- Let's think of a function that converts $\mathbb{R} \mapsto [0, 1]$
- **Sigmoid function** looks nice: $\sigma(z) = \frac{1}{1 + e^{-z}}$
- Wrap it around: $a_\sigma(\mathbf{x}, \boldsymbol{\theta}) = \sigma(a(\mathbf{x}, \boldsymbol{\theta})) = \frac{1}{1 + e^{-\langle \mathbf{x}, \boldsymbol{\theta} \rangle}}$
- Now it is possible to plug a_σ into log-likelihood and train:
$$\hat{a} = \arg \max_{\boldsymbol{\theta} \in \Theta} \log \mathcal{L}(\boldsymbol{\theta}) =$$
$$= \arg \max_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^N y_i \log \sigma(\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle) + (1 - y_i) \log (1 - \sigma(\langle \mathbf{x}_i, \boldsymbol{\theta} \rangle))$$
- This approach is called **logistic regression**



How to make classifier return probabilities?

- $a(x, \theta)$ has a range of \mathbb{R} , but it should be $[0, 1]$
- Let's think of a function that converts $\mathbb{R} \mapsto [0, 1]$
- **Sigmoid function** looks nice: $\sigma(z) = \frac{1}{1 + e^{-z}}$
- Wrap it around: $a_\sigma(x, \theta) = \sigma(a(x, \theta)) = \frac{1}{1 + e^{-\langle x, \theta \rangle}}$
- Now it is possible to plug a_σ into log-likelihood and train:
$$\hat{a} = \arg \max_{\theta \in \Theta} \log \mathcal{L}(\theta) =$$
$$= \arg \max_{\theta \in \Theta} \sum_{i=1}^N y_i \log \sigma(\langle x_i, \theta \rangle) + (1 - y_i) \log (1 - \sigma(\langle x_i, \theta \rangle))$$
- This approach is called **logistic regression**



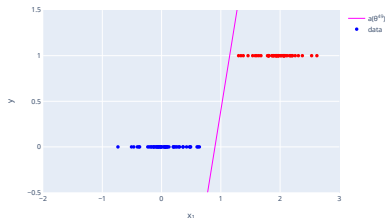
NB: usage of σ and log-likelihood in general does not guarantee probabilistic interpretation

Comparison

Linear classifier

$$Q(\theta) = \sum_{i=1}^N [\text{sign}(a(x_i, \theta)) \neq y_i]$$

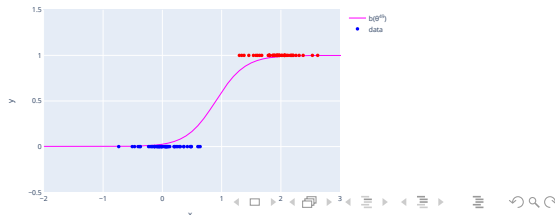
- $a(x, \theta) = \langle x, \theta \rangle$
- $a(x, \theta) : \mathbb{R}^K \times \mathbb{R}^K \mapsto \mathbb{R}$
- ✓ Easy to interpret
- ✗ Gives only yes-no answer



Logistic regression

$$\log \mathcal{L}(\theta) = \sum_{i=1}^N y_i \log a_\sigma(x_i, \theta) + (1 - y_i) \log (1 - a_\sigma(x_i, \theta))$$

- $a_\sigma(x, \theta) = \sigma[a(x, \theta)] = \sigma[\langle x, \theta \rangle]$
- $a_\sigma(x, \theta) : \mathbb{R}^K \times \mathbb{R}^K \mapsto [0, 1]$
- ✓ More profound interpretation
- ✓ Returns probability of positive class $y = 1$

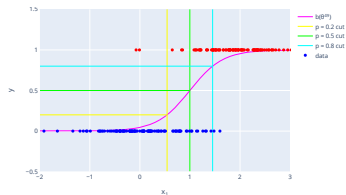


Output of binary classifier

Let **red** correspond to a positive class (1)
and **blue** – to a negative class (0)

- 1 classified as 1 \Rightarrow True positive (TP)
- 0 classified as 0 \Rightarrow True negative (TN)
- 1 classified as 0 \Rightarrow False negative (FN)
- 0 classified as 1 \Rightarrow False positive (FP)

NB: in case of logistic regression you need to specify a **cut** on the model's output to get deterministic prediction



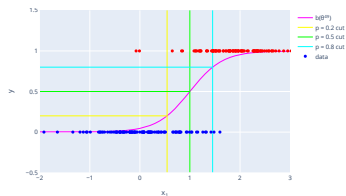
Output of binary classifier

Let **red** correspond to a positive class (1)
and **blue** – to a negative class (0)

- 1 classified as 1 \Rightarrow True positive (TP)
- 0 classified as 0 \Rightarrow True negative (TN)
- 1 classified as 0 \Rightarrow False negative (FN)
- 0 classified as 1 \Rightarrow False positive (FP)

NB: in case of logistic regression you need to specify a **cut** on the model's output to get deterministic prediction

→ Then how do we measure model's performance?



Metrics

Using these output for binary classifier we can construct several quality metrics:

- **accuracy** = $\frac{TP + TN}{TP + TN + FP + FN}$
→ fraction of correct answers
- **precision** = $\frac{TP}{TP + FP}$
→ fraction of correct positive answers in all positive answers
- **recall** = $\frac{TP}{TP + FN}$
→ fraction of correctly classified positive objects
- **F1-score** = $\frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$
→ harmonic mean of precision and recall



Metrics

Using these output for binary classifier we can construct several quality metrics:

- **accuracy** = $\frac{TP + TN}{TP + TN + FP + FN}$
→ fraction of correct answers
- **precision** = $\frac{TP}{TP + FP}$
→ fraction of correct positive answers in all positive answers
- **recall** = $\frac{TP}{TP + FN}$
→ fraction of correctly classified positive objects
- **F1-score** = $\frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$
→ harmonic mean of precision and recall



choose metrics wisely!

Example —○ what if dataset is imbalanced?

- Consider a dataset with 10000 **negative** and 10 **positive** objects

Example —○ what if dataset is imbalanced?

- Consider a dataset with 10000 **negative** and 10 **positive** objects
- After training you might end up with a model which returns **negative** answers all the time

Example —○ what if dataset is imbalanced?

- Consider a dataset with 10000 **negative** and 10 **positive** objects
- After training you might end up with a model which returns **negative** answers all the time
- Because from loss perspective it might be optimal – contribution of misclassified positives to the loss function will be marginal:

$$\log \mathcal{L}(\boldsymbol{\theta}) \sim \sum_{i=1}^{10} \log a_{\sigma}(\mathbf{x}_i, \boldsymbol{\theta}) + \sum_{i=1}^{10000} \log (1 - a_{\sigma}(\mathbf{x}_i, \boldsymbol{\theta}))$$

Example —○ what if dataset is imbalanced?

- Consider a dataset with 10000 **negative** and 10 **positive** objects
- After training you might end up with a model which returns **negative** answers all the time
- Because from loss perspective it might be optimal – contribution of misclassified positives to the loss function will be marginal:

$$\log \mathcal{L}(\boldsymbol{\theta}) \sim \sum_{i=1}^{10} \log a_{\sigma}(\mathbf{x}_i, \boldsymbol{\theta}) + \sum_{i=1}^{10000} \log (1 - a_{\sigma}(\mathbf{x}_i, \boldsymbol{\theta}))$$

- Furthermore, if you used accuracy as a metric to be optimised, you would get:

$$\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} = \frac{0 + 10000}{0 + 10000 + 0 + 10} \approx \mathbf{99.9\%!}$$

Example —○ what if dataset is imbalanced?

- Consider a dataset with 10000 **negative** and 10 **positive** objects
- After training you might end up with a model which returns **negative** answers all the time
- Because from loss perspective it might be optimal – contribution of misclassified positives to the loss function will be marginal:

$$\log \mathcal{L}(\boldsymbol{\theta}) \sim \sum_{i=1}^{10} \log a_{\sigma}(\mathbf{x}_i, \boldsymbol{\theta}) + \sum_{i=1}^{10000} \log (1 - a_{\sigma}(\mathbf{x}_i, \boldsymbol{\theta}))$$

- Furthermore, if you used accuracy as a metric to be optimised, you would get:

$$\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} = \frac{0 + 10000}{0 + 10000 + 0 + 10} \approx \mathbf{99.9\%!}$$

- Hence, you could've assumed that the model is really good (but it's not)

Example —○ what if dataset is imbalanced?

more in homework!

- Consider a dataset with 10000 **negative** and 10 **positive** objects
- After training you might end up with a model which returns **negative** answers all the time
- Because from loss perspective it might be optimal – contribution of misclassified positives to the loss function will be marginal:

$$\log \mathcal{L}(\boldsymbol{\theta}) \sim \sum_{i=1}^{10} \log a_{\sigma}(\mathbf{x}_i, \boldsymbol{\theta}) + \sum_{i=1}^{10000} \log (1 - a_{\sigma}(\mathbf{x}_i, \boldsymbol{\theta}))$$

- Furthermore, if you used accuracy as a metric to be optimised, you would get:

$$\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} = \frac{0 + 10000}{0 + 10000 + 0 + 10} \approx \mathbf{99.9\%!}$$

- Hence, you could've assumed that the model is really good (but it's not)
- Always take imbalance into account! (e.g. up-/downsampling, reweighting)
- Be careful when interpreting your metrics!

Choosing the right threshold

- As was mentioned before, classifier decision (thus performance) **depends on the threshold**
- How can we estimate its performance *for all cuts* in general?
 - define **true positive rate**: $TPR = \frac{TP}{TP + FN}$
 - and **false positive rate**: $FPR = \frac{FP}{FP + TN}$
 - move the **green** line on the plot
 - moving to the left results in $TP \uparrow$ and $FN \downarrow$ (sensitive, but not specific)
 - moving to the right results in $TP \downarrow$ and $FN \uparrow$ (specific, but not sensitive)



Choosing the right threshold

- As was mentioned before, classifier decision (thus performance) **depends on the threshold**
- How can we estimate its performance *for all cuts* in general?

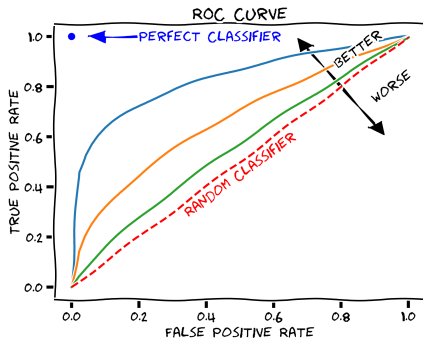
- define **true positive rate**: $TPR = \frac{TP}{TP + FN}$
- and **false positive rate**: $FPR = \frac{FP}{FP + TN}$
- move the **green** line on the plot
- moving to the left results in $TP \uparrow$ and $FN \downarrow$ (sensitive, but not specific)
- moving to the right results in $TP \downarrow$ and $FN \uparrow$ (specific, but not sensitive)



Idea: plot jointly values of TPR and FPR for different values of threshold

Receiver Operating Characteristic — o curve

- Scanning through thresholds we can calculate for each of them TPR and FPR and plot them
 - This will give a **ROC curve**
- Random classifier – diagonal line at ROC
- Perfect classifier – step function
- Something-went-wrong classifier – curve below diagonal line

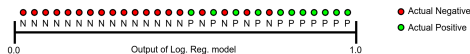


Receiver Operating Characteristic —○ area under curve

The power of a classifier can be represented by the **area under ROC curve**:

$$\text{ROC AUC} = \int_0^1 \text{TPR}(\text{FPR}) d(\text{FPR})$$

- aggregated measure of performance across all possible thresholds \Rightarrow **threshold-invariant**
- 1 \Rightarrow predictions are 100% correct
- 0.5 \Rightarrow random classifier
- 0 \Rightarrow predictions are 100% wrong
- probability that a random positive (**green**) example is positioned to the right of a random negative (**red**) example \Rightarrow **scale-invariant**

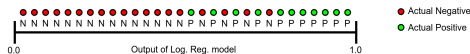


Receiver Operating Characteristic — area under curve

The power of a classifier can be represented by the **area under ROC curve**:

$$\text{ROC AUC} = \int_0^1 \text{TPR}(\text{FPR}) d(\text{FPR})$$

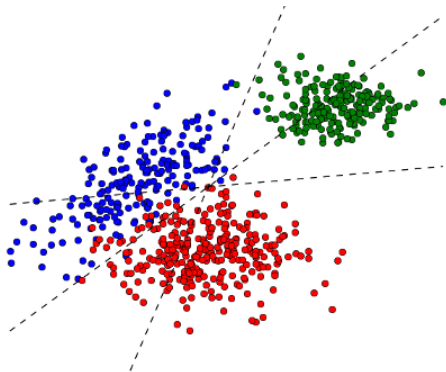
- aggregated measure of performance across all possible thresholds \Rightarrow **threshold-invariant**
- 1 \Rightarrow predictions are 100% correct
- 0.5 \Rightarrow random classifier
- 0 \Rightarrow predictions are 100% wrong
- probability that a random positive (**green**) example is positioned to the right of a random negative (**red**) example \Rightarrow **scale-invariant**



Nice moment to show this animation

Multiclass classification

- If the number of classes > 2 , then we need to modify classification approach
- $|Y| < \infty = \{1, \dots, L\} \Rightarrow$ impossible to use single model



From multiclass problem to a set of binary problems

"One vs all"

- Train L binary classifiers:
 $b_\ell(\mathbf{x}, \boldsymbol{\theta}), \ell \in \{1, \dots, L\}$
- Each $a_\ell(\mathbf{x}, \boldsymbol{\theta})$ is trained to separate between ℓ -th class and all the rest
- The most confident classifier gives the final separation:

$$a(\mathbf{x}, \boldsymbol{\theta}) = \arg \max_{\ell \in \{1, \dots, L\}} b_\ell(\mathbf{x}, \boldsymbol{\theta})$$

"All vs all"

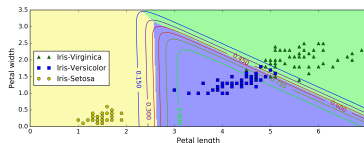
- Train C_L^2 binary classifiers to separate every two classes:
 $a_{i,j}(\mathbf{x}, \boldsymbol{\theta}), i, j \in 1, \dots, L, i \neq j$
- In this case each classifier decides only between two classes
- $\hat{y} = \arg \max_{\ell \in \{1, \dots, L\}} \sum_{i=1}^L \sum_{j \neq i} [b_{i,j}(\mathbf{x}, \boldsymbol{\theta}) = \ell]$
- Class with the highest vote over C_L^2 classifiers wins!

Multiclass logistic regression

- Suppose we trained $a_\ell(\mathbf{x}, \boldsymbol{\theta}) = \langle \mathbf{x}, \boldsymbol{\theta} \rangle$, $\ell = 1, \dots, L$ each separating its individual class
- Now we need to turn their output into probabilities
- Reminder: $a_\ell(\mathbf{x}, \boldsymbol{\theta})$ has a range \mathbb{R}
- We need to convert vector $[a_1(\mathbf{x}, \boldsymbol{\theta}), \dots, a_L(\mathbf{x}, \boldsymbol{\theta})]$ into a vector of probabilities
- We can put **softmax** on top:

$$[a_1(\mathbf{x}, \boldsymbol{\theta}), \dots, a_L(\mathbf{x}, \boldsymbol{\theta})] \mapsto \left[\frac{\exp[a_1(\mathbf{x}, \boldsymbol{\theta})]}{\sum_{\ell=1}^L \exp[a_\ell(\mathbf{x}, \boldsymbol{\theta})]}, \dots, \frac{\exp[a_L(\mathbf{x}, \boldsymbol{\theta})]}{\sum_{\ell=1}^L \exp[a_\ell(\mathbf{x}, \boldsymbol{\theta})]} \right]$$

- Then components sum up to 1 and each component represents* probability that object \mathbf{x} belongs to a correspondent class

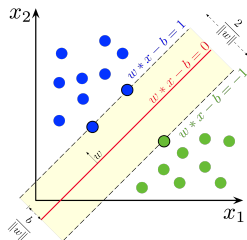
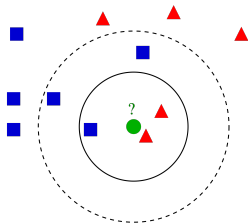


Final note

We didn't have time to describe them, but in the family of "classical algorithms" there are also these guys – and they are cool enough to be studied and tried out:

- **kNN** – k -nearest neighbors algorithm
- **SVM** – Support Vector Machine
- **Naive Bayes** – Naive Bayes classifier

$$\hat{y} = \arg \max_y P(y) \prod_{i=1}^n P(x_i | y)$$



Summary

- **Linear models** are the simplest example of models used for **SL**
- Solution of optimization problem can be found analytically, but **(stochastic) gradient descent** is a more general approach to training
- During the process of training model's performance is checked **throughout** to control **under-** and **overtraining**
- Model **regularization** prevents overtraining thanks to additional term in the loss function
- Several variations to MSE and MAE loss functions exist, each with its unique properties – all can be used both **as losses and as metrics**

Summary

- **Linear models** are the simplest example of models used for **SL**
- Solution of optimization problem can be found analytically, but **(stochastic) gradient descent** is a more general approach to training
- During the process of training model's performance is checked **throughout** to control **under-** and **overtraining**
- Model **regularization** prevents overtraining thanks to additional term in the loss function
- Several variations to MSE and MAE loss functions exist, each with its unique properties – all can be used both **as losses and as metrics**
- Classification problem has pretty the same pipeline as the regression. Furthermore, it can be converted into **logistic regression** that estimates **probability** of belonging to a class
- **Dedicated classification metrics** exist, all being specific/sensitive to its own thing, both threshold dependent (e.g. accuracy) and not (e.g. ROC AUC)
- **Multiclass clasification** can be also performed using "one-vs-all" and "all-vs-all" methods, or in case of log.regression with applying **softmax function**