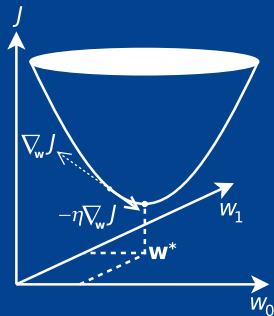


# MACHINE LEARNING

## LESSON 5: Training II

CARSTEN EIE FRIGAARD  
SPRING 2019



# L05: Training II

## Agenda

- ▶ Important: your feedback (positive as well as negative)...
- ▶ Course modification i): *more OPTIONAL* exercises
- ▶ Course modification ii): revision to 'Microlearning'  
(only one lecture session, beginning of class)
- ▶ Installing keras is slow/buggy, see WIKI..
- ▶ L04 Training I: Training a linear regression model  
OPTIONAL: Exercise: [L04/linear\\_regression\\_1.ipynb](#)  
OPTIONAL: Exercise: [L04/linear\\_regression\\_2.ipynb](#)
- ▶ L05 Training II: **Training and model concepts**  
Exercise: [L05/gradient\\_descent.ipynb](#)  
Exercise: [L05/capacity\\_under\\_overfitting.ipynb](#)  
Exercise: [L05/generalization\\_error.ipynb](#)  
OPTIONAL: Exercise: [L05/train\\_test\\_split.ipynb](#)

# RESUMÉ: Metrics

Precision, recall and accuracy,  $F_1$ -score, and confusion matrix

$$\text{precision, } p = \frac{TP}{TP+FP}$$

$$\text{recall (or sensitivity), } r = \frac{TP}{TP+FN}$$

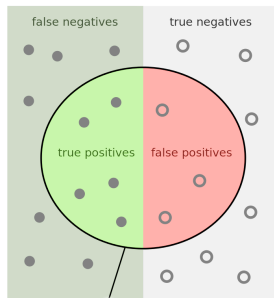
$$\text{accuracy, } a = \frac{TP+TN}{TP+TN+FP+FN}$$

$$F_1\text{-score, } F_1 = \frac{2pr}{p+r}$$

Confusion Matrix, binary-class data,

$M_{\text{confusion}} =$

	actual true	actual false
predicted true	TP	FP
predicted false	FN	TN



$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

# RESUMÉ: Covariance Matrix

Data matrix for a two-dimensional feature space

$$\mathbf{X} = \begin{bmatrix} \overset{\lambda_1}{x_1^{(1)}} & \overset{\lambda_2}{x_2^{(1)}} \\ x_1^{(2)} & x_2^{(2)} \\ \vdots & \vdots \\ x_1^{(n)} & x_2^{(n)} \end{bmatrix} = \begin{bmatrix} 0.35 & -7.62 \\ -4.99 & 13.79 \\ \vdots & \vdots \\ 9.54 & -25.64 \\ 4.21 & -2.25 \end{bmatrix}$$

Covariance matrix, for the two-dimensional feature space

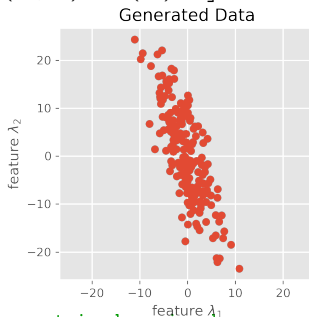
$$\mathbf{\Sigma}(\mathbf{X}) = \begin{bmatrix} \sigma(\lambda_1, \lambda_1) & \sigma(\lambda_1, \lambda_2) \\ \sigma(\lambda_2, \lambda_1) & \sigma(\lambda_2, \lambda_2) \end{bmatrix} = \begin{bmatrix} \sigma(\lambda_1)^2 & \sigma(\lambda_1, \lambda_2) \\ \sigma(\lambda_2, \lambda_1) & \sigma(\lambda_2)^2 \end{bmatrix}$$

with 
$$\sigma(\lambda_1, \lambda_2) = \frac{1}{n} \sum_{i=1}^n (x_1^{(i)} - \mu_{\lambda_1})(x_2^{(i)} - \mu_{\lambda_2})$$

Example:  $\mathbf{X}$ ; a 100 x 2 matrix, see fig..

$$\mathbf{\Sigma}(\mathbf{X}) = \begin{bmatrix} 13.2 & -28.8 \\ -28.8 & 93.3 \end{bmatrix}$$

- ▶  $\mathbf{\Sigma}$  is real and symmetric,
- ▶ diagonal: the (auto)-variance of a feature,  $\sigma(\lambda_j)^2$
- ▶ Pearson's  $r$ : cross-correlation via cross-covar,
- ▶ ~~similar dimension as Confusion matrix,~~
- ▶ python implementation: see [L02/Extra/covariance\\_matrix\\_demo.ipynb](#).



# RESUMÉ: Covariance Matrix, Take II

For a dataset,  $\mathbf{X}$ ; features cat and dog; classifying cat/non-cat

Covariance matrix (dim = features x features)

$$\Sigma(\mathbf{X}) = \begin{bmatrix} \sigma_{\text{cat,cat}} & \sigma_{\text{cat,dog}} \\ \sigma_{\text{dog,cat}} & \sigma_{\text{dog,dog}} \end{bmatrix} = \begin{array}{c|cc} & \text{cat} & \text{dog} \\ \hline \text{cat} & \sigma_{\text{cat}}^2 & \sigma_{\text{cat,dog}} \\ \text{dog} & \sigma_{\text{dog,cat}} & \sigma_{\text{dog}}^2 \end{array}$$

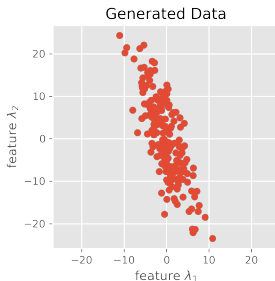
co-variance

$$\sigma_{\text{cat,dog}} = \frac{1}{n} \sum_{i=1}^n (x_{\text{cat}}^{(i)} - \mu_{\text{cat}})(x_{\text{dog}}^{(i)} - \mu_{\text{dog}})$$

Example  $\Sigma(\mathbf{X}) = \begin{bmatrix} 13.2 & -28.8 \\ -28.8 & 93.3 \end{bmatrix}$

Confusion matrix, cat/non-cat  
(dim = classes x classes)

$$\mathbf{M} = \begin{bmatrix} \text{TP} & \text{FP} \\ \text{FN} & \text{TN} \end{bmatrix} = \begin{array}{c|cc} & \text{actual} & \\ \hline & \text{cat} & \text{non-cat} \\ \hline \text{cat} & \text{T, cat} & \text{F, cat} \\ \text{dog} & \text{F, dog} & \text{T, non-cat} \end{array}$$



# RESUMÉ: Covariance Matrix

- ▶ python implementation: see [L02/Extra/covariance\\_matrix\\_demo.ipynb](#).

not so scary, afterall...

```
1  # Covariance
2  def cov(x, y, bias=True):
3      assert len(x)==len(y)
4      xbar, ybar = x.mean(), y.mean()
5      if bias:
6          n=len(x)
7      else:
8          n=len(x) - 1
9      assert n>0
10     return np.sum((x - xbar)*(y - ybar))/n
11
12 # Covariance matrix
13 def cov_mat(X, bias=True):
14     return np.array([ \
15         [cov(X[0], X[0], bias), cov(X[0], X[1], bias)], \
16         [cov(X[1], X[0], bias), cov(X[1], X[1], bias)] \
17     ])
```

# RESUMÉ: Training a Linear Regressor

## Minimizing the Linear Regression: The `argmin` concept

Our linear regression cost function was

$$J(\mathbf{X}, \mathbf{y}; \mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

and training amounts to finding a value of  $\mathbf{w}$ , that minimizes  $J$ . This is denoted as

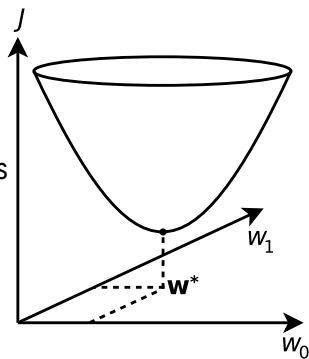
$$\begin{aligned}\mathbf{w}^* &= \operatorname{argmin}_{\mathbf{w}} J(\mathbf{X}, \mathbf{y}; \mathbf{w}) \\ &= \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2\end{aligned}$$

and by minima, we naturally hope for

► the global minimum

thought for non-linear models this cannot be guaranteed, hitting some

► local minimum



# RESUMÉ: L04/linear\_regression\_1.ipynb

## Training: The Closed-form Linear-Least-Squares Solution

To solve for  $\mathbf{w}^*$  in closed form, we find the gradient of  $J$  with respect to  $\mathbf{w}$

$$\nabla_{\mathbf{w}} J = \left[ \frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \dots, \frac{\partial J}{\partial w_m} \right]^\top$$

Taking the partial derivery  $\partial/\partial_{\mathbf{w}}$  of the  $J$  via the gradient (nabla) operator

$$\begin{aligned} \nabla_{\mathbf{w}} J(\mathbf{X}, \mathbf{y}; \mathbf{w}) &= \mathbf{X}^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) = 0 \\ 0 &= \mathbf{X}^\top \mathbf{X}\mathbf{w} - \mathbf{X}^\top \mathbf{y} \end{aligned}$$

with a small amount of matrix algebra, this gives the *normal equation*

$$\begin{aligned} \mathbf{w}^* &= \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \end{aligned}$$



# Numerical Solution: Gradient Descent (GD)

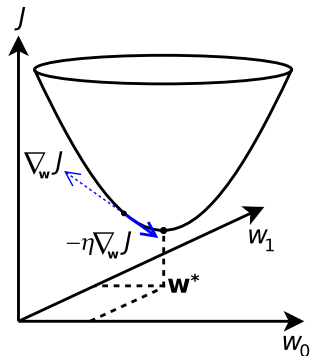
## The GD Algorithm

First, find the deriverty of  $J$ , via  $\nabla_{\mathbf{w}}J$ , that after some matrix algebra gives

$$\nabla_{\mathbf{w}}J(\mathbf{w}) = \frac{2}{m}\mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{y})$$

Then move along in the *opposite* direction of this gradient, taking a step of size  $\eta$

$$\mathbf{w}^{(\text{step } N+1)} = \mathbf{w}^{(\text{step } N)} - \eta \nabla_{\mathbf{w}}J(\mathbf{w})$$



The SG-algo in python code

```
1 for iteration in range(n_iterations):  
2     gradients=2/m*X_b.T.dot(X_b.dot(theta)-y)  
3     theta=theta - eta * gradients
```

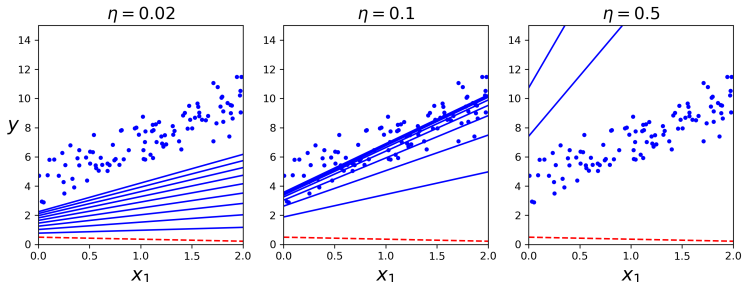
**NOTE:** The  $\mathbf{X\_b}$  is a  $\mathbf{X}$  with a all-1 column prepended, and using the contant factor  $2/m$  instead of just  $1/2$  and with  $\mathbf{w} = \boldsymbol{\theta}$

# Learning rate, $\eta$

And understanding Scikit-learn's SGDRegressor

Scikit-learn  $\eta$  updating schemes:

**constant**, adaptive, invscaling, optimal



The SGDRegressor in Scikit-learn has  
a *hyperparameter* for this:

`learning_rate: string, optional`

`'invscaling': [default], eta=eta0/pow(t,power_t)`



# Learning rate, $\eta$

And understanding Scikit-learn's SGDRegressor

The SGDRegressor constructor in Scikit-learn

```
1 class sklearn.linear_model.SGDRegressor(  
2     loss      = 'squared_loss', penalty      = 'l2',  
3     alpha    = 0.0001,      l1_ratio      = 0.15,  
4     tol      = None,        shuffle      = True,  
5     verbose  = 0,           epsilon      = 0.1,  
6     eta0     = 0.01,        power_t     = 0.25,  
7     n_iter_no_change=5,     warm_start  = False,  
8     fit_intercept = True,   max_iter    = None,  
9     average     = False,   n_iter      = None  
10    random_state  = None,   learning_rate='invscaling',  
11    early_stopping=False,  validation_fraction=0.1  
12 )
```



Important for now...(hyperparam search in L07)

- ▶ loss, penalty (our MSE and  $\mathcal{L}_2$  norm),
- ▶ eta0, learning\_rate,
- ▶ shuffle, early\_stopping,
- ▶ and perhaps random\_state.

# Stochastic Gradient Descent (SGD) Method

Exercise: `gradient_descent.ipynb`

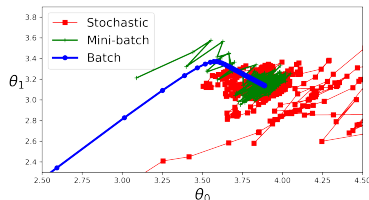
The problem with the GD Algorithm, it takes  $\mathbf{X}$  as input, the complete dataset

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \frac{2}{m} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

That's a big mouthful to matrix transpose and multiply!  
Introducing a per-data-sample, or stochastic, method

$$\nabla_{\mathbf{w}} J(\mathbf{w})^{(i)} = \frac{2}{m} \mathbf{x}^{(i)T} (\mathbf{x}^{(i)}\mathbf{w} - y)$$

```
1 for epoch in range(n_epochs):
2     for i in range(m):
3         ...
4         r = np.random.randint(m)
5         xi = X_b[r:r+1]
6         yi = y[r:r+1]
7         grads = 2*xi.T.dot(xi.dot(theta)-yi)
8         eta = ...
9         theta = ...
```



**NOTE:** Notice the use of `epoch` in SGD, that is different from `iteration` in GD.

# Model capacity and under/overfitting

Exercise: `capacity_under_overfitting.ipynb`

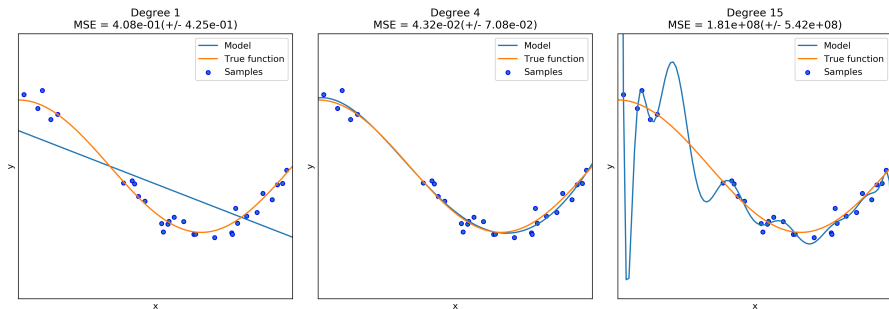
Dummy and Paradox classifier:

*capacity fixed  $\sim 0$ , cannot generalize at all!*

Linear regression for a polynomial model:

*capacity  $\sim$  degree of polynomial,  $x^n$*

Polynomial linear reg. fit for underlying model:  $\cos(x)$

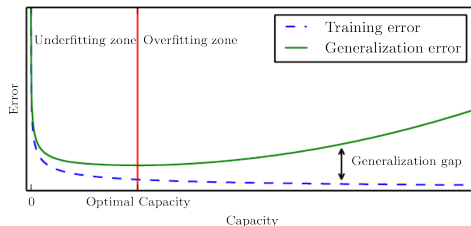
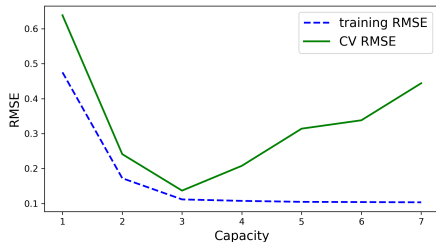


- ▶ underfitting: capacity of model too low,
- ▶ overfitting: capacity too high.

# Generalization Error

Exercise: `generalization_error.ipynb`

RMSE-capacity plot for lin. reg. with polynomial features  
(capacity = degree of poly)



Inspecting the plots from [HOML] and [DL], extracting the concepts

- ▶ training/generalization error,
- ▶ generalization gap,
- ▶ underfit/overfit zone,
- ▶ optimal capacity (best-model, early stop).

