regularized regression: we want to minimize the test error (generalization error).
Usually, the test error consists of several contributions, for example
bias and variance.

observation: the biggest error source dominates the total ("a chain is
only as strong as it's weakest link"). It makes no sense to
make all errors except for one small - wash of resources.
All error sources should contribute about equally.

In OLS: bias $= 0$, variance can be high if features are (nearly)
redundant (i.e. linearly dependent).

$\Rightarrow$ allow some bias if this reduces the variance a lot.

regularization : add a new term to the loss fct. that reduces overfitting

original Loss = data term , | new Loss = data term + $\tau \cdot$ regularization term
                            how well does the model                    prior knowledge: how
                            with parameters $\beta$ fit the training set ?    should models for this problem
                                                                             class typically look like ?
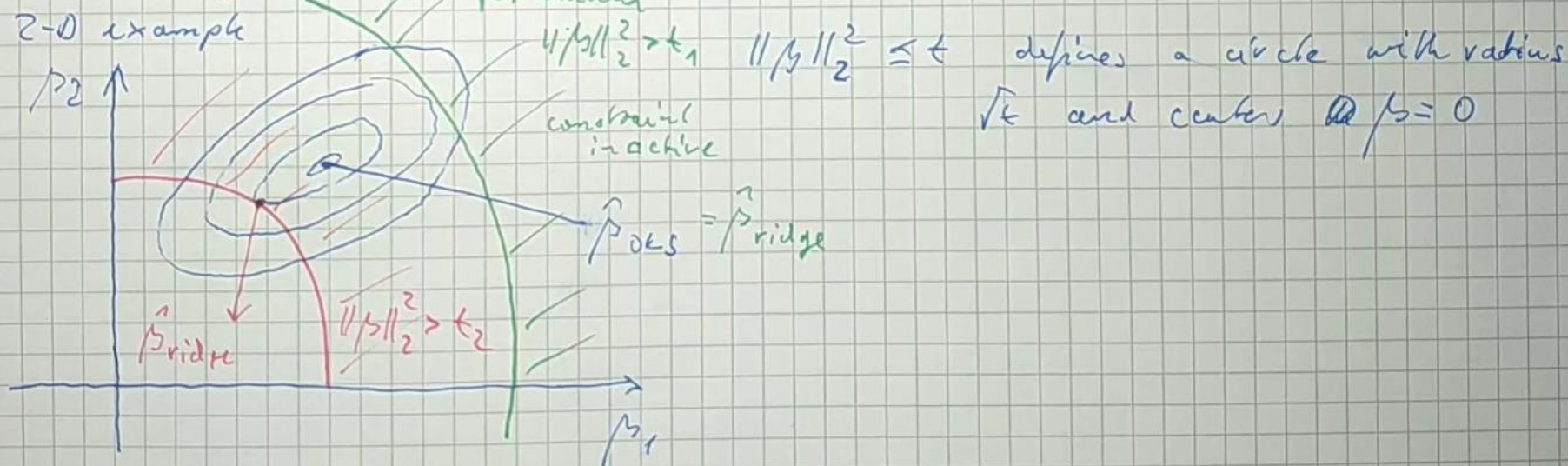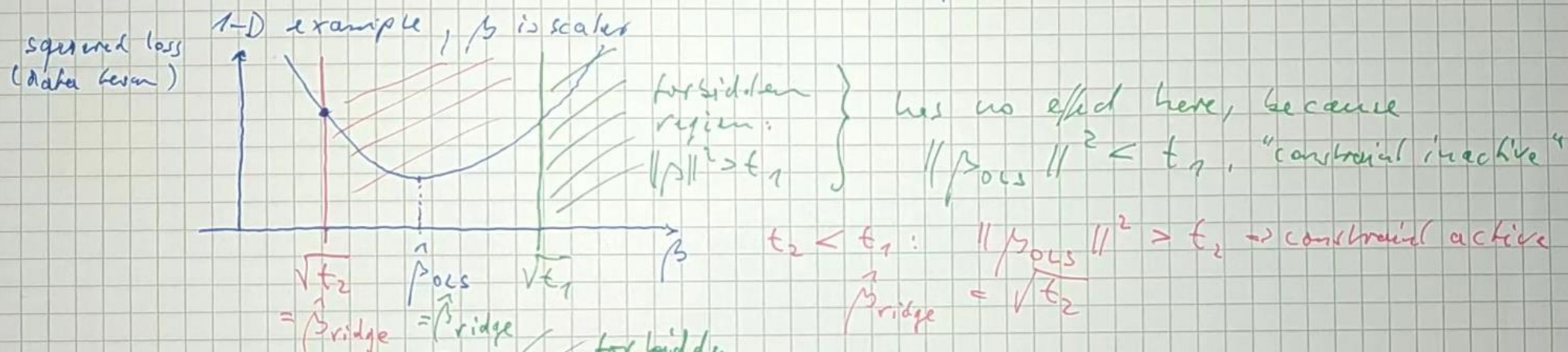                                                                             which parameters $\beta$ are plausible?

ridge regression: data term $\hat{=}$ squared loss $\|Y - X\beta\|^2$

regularization term $\hat{=}$ $L_2$ regularizer · $\|\beta\|_2^2 = \beta^T \beta$

interpretation in term of prior knowledge: - data term $\hat{=}$ negative log-likelihood
of the TS $\Rightarrow$ regularization is the neg. log likelihood of $\exp(-\tau \beta^T \beta)$
$\Rightarrow$ prior belief: good parameters are Gaussian distributed with mean $0$ and variance $\sigma^2 = \frac{1}{2\tau}$

graphical illustration of ridge reg. $\hat{\beta} = \arg\min_{\beta} \|Y - X\beta\|^2 + \tau \|\beta\|_2^2$

equivalent to $\hat{\beta} = \arg\min_{\beta} \|Y - X\beta\|^2$ s.t. $\|\beta\|_2^2 \le t$

1-D example, $\beta$ is scalar

squared loss (data term)



forbidden region: $\|\beta\|^2 > t_1$ } has no effect here, because $\|\hat{\beta}_{OLS}\|^2 < t_1$, "constraint inactive"

$t_2 < t_1$: $\|\hat{\beta}_{OLS}\|^2 > t_2 \Rightarrow$ constraint active

$\hat{\beta}_{ridge} = \sqrt{t_2}$

$\sqrt{t_2}$   $\hat{\beta}_{OLS}$   $\sqrt{t_1}$   $\beta$
$= \hat{\beta}_{ridge}$   $= \hat{\beta}_{ridge}$

2-D example

$\beta_2$



forbidden $\|\beta\|_2^2 > t_1$

constraint inactive

$\hat{\beta}_{OLS} = \hat{\beta}_{ridge}$

$\hat{\beta}_{ridge}$

$\|\beta\|_2^2 > t_2$

$\beta_1$

$\|\beta\|_2^2 \le t$ defines a circle with radius $\sqrt{t}$ and center @ $\beta = 0$

other possibilities for regularization: feature selection, $L_1$ regularization

**feature selection:** idea: when features are redundant, only use the most relevant
non-redundant subset = "active set" $A \subseteq \{1, ..., D\}$

objective: $\hat{\beta} = \arg\min_{\beta} \|Y - X\beta\|^2$ s.t. $|A| \leq t$

$\underbrace{\qquad}_{\text{number of active features}}$

(only $\beta_j \neq 0$ when $j \in A$, otherwise $\underbrace{\beta_j = 0}_{\text{inactive coefficients}}$)

"0-norm": $\|\beta\|_0 = \#$ non-zero coefficients in $\beta = |A|$

$\hat{\beta} = \arg\min_{\beta} \|Y - X\beta\|^2 + \tau \|\beta_0\|$

advantage: $\hat{\beta}$ is the exact OLS solution for the active features (= unbiased)

disadvantage: $L_0$ regularization is NP-hard, $\Rightarrow$ the active set for
some $t$ can be very different than for $t-1$ $\Rightarrow$ in general, finding
the globally optimal active set for some $t$ requires **exhaustive search**
in practice, things are usually not so bad: optimal active sets for $t-1$
and $t$ are usually very similar $\Rightarrow$ efficient approximation alg.

**orthogonal matching pursuit (OMP):** ⓪ start with $A = \emptyset$     OLS solution for $A^{(m-1)}$

① for $m = 1, ..., t$: add one feature per iteration:
ⓐ compute ~~squared~~ residuals for current guess: $R_i = \left(Y_i - X_i \beta^{(m-1)}\right)$
ⓑ find the best inactive feature: $\hat{j} = \arg\max_{j \notin A^{(m-1)}} |X_{\hat{j}}^T R|$ and add: $A^{(m)} = A^{(m-1)} \cup \{\hat{j}\}$

$L_1$ regularization     LASSO regression   ("least absolute shrinkage and selection operator")

in between   $L_2$ regularization $\|\beta\|_2^2 \leq t$   and $L_0$ regul. $\|\beta\|_0 \leq t$

it shrinks coefficients towards 0 like $L_2$ reg. , but not as much

selects features like $L_0$ reg. , but is not NP-hard , but convex

objective :   $\hat{\beta} = \underset{\beta}{\arg\min} \|Y - X\beta\|^2$   s.t.  $\|\beta\|_1 \leq t$

$\hat{\beta} = \underset{\beta}{\arg\min} \|Y - X\beta\|^2 + \tau \|\beta\|_1$     $\|\beta\|_1 = \sum_{j=1}^{D} |\beta_j|$

convex $\Rightarrow$ unique solution , easy to find , various algo- , for example

LARS algorithm : similar to OMP , but after each iteration , check if $A^{(m)}$ has become redundant after adding the new feature ( in practice , this occurs every once in a while , but not so very $O(\log n)$ )

if yes $\Rightarrow$ remove the least important feature from $A^{(m)}$ before adding a new one

in 1-D , LASSO and ridge regression are identical, because $\|\beta\|_2 = \|\beta\|_1$

$\Rightarrow$ same diagram

allowed region $\|\beta\|_1 \leq t$ : diamond

$\hat{\beta}_{OLS} = \hat{\beta}_{LASSO}$

inactive constraint

$\|\beta\|_1 \leq t_1$

$L_2$ reg.    $\|\beta\|_1 \leq t_2$

$\hat{\beta}_{LASSO}$ $(\beta_1 = 0)$ $\Rightarrow$ feature selection

$\sqrt{t_2}$  $\beta_{OLS}$  $\sqrt{t_1}$  $\beta$