Example and Task
Idealized regression
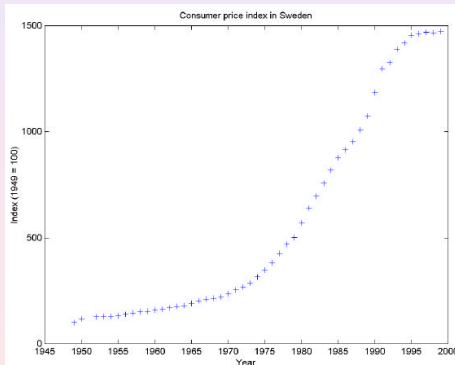Error measures
The real regression
Examples
Model selection

# Introduction to Regression

Antanas Verikas
antanas.verikas@hh.se

IDE, Halmstad University

2013

Example and Task
Idealized regression
Error measures
The real regression
Examples
Model selection

Example
Task
Data and assumptions

## Example

Regression aims at finding a function that fits the observations.



Figure: Consumer prise index in Sweden.

Observations:
(x,y) pairs

(1949, 100)
(1950, 117)
...
(1996, 1462)
(1997, 1469)
(1998, 1467)
(1999, 1474)

Example and Task
Idealized regression
Error measures
The real regression
Examples
Model selection

Example
Task
Data and assumptions

## Linear fit

The linear fit is not so good.



| $y$ | $\widehat{y}$ |
|------|------|
| 100 | -215 |
| 117 | -184 |
| ... | ... |
| 1467 | 1314 |
| 1474 | 1345 |

Figure: Consumer prise index in Sweden, linear fit.

Example and Task
Idealized regression
Error measures
The real regression
Examples
Model selection

Example
Task
Data and assumptions

## Example

Apply a transformation.



Take logarithm of $y$ and fit a straight line.

Figure: Consumer prise index in Sweden.

Example and Task
Idealized regression
Error measures
The real regression
Examples
Model selection

Example
Task
Data and assumptions

## Linear fit

Transform $y$ back to the original. The fit is better.



| $y$ | $\widehat{y}$ |
|------|------|
| 100 | 83 |
| 117 | 88 |
| ... | ... |
| 1467 | 1660 |
| 1474 | 1765 |

Figure: Consumer prise index in Sweden.

Example and Task
Idealized regression
Error measures
The real regression
Examples
Model selection

Example
Task
Data and assumptions

## Regression task

Construct a model of a process, using examples of the process.

Input: $\mathbf{x}$ (possibly a vector)

Output: $y = g(\mathbf{x})$ (generated by the process)

Examples: Pairs of input and output $\{y(n), \mathbf{x}(n)\}$

Our model: $\widehat{y} = f(\mathbf{x})$

The function $f$ is our estimate of the true function $g$

Example and Task
Idealized regression
Error measures
The real regression
Examples
Model selection

Example
Task
Data and assumptions

## Data and assumptions

$$\text{Data set } \mathbf{Z} = \{\mathbf{x}(n), y(n)\}_{n=1,\ldots,N}$$
$$y(n) = g[(\mathbf{x}(n)] + \varepsilon(n)$$

$\mathbf{x}(n)$  Observed input

$y(n)$  Observed output

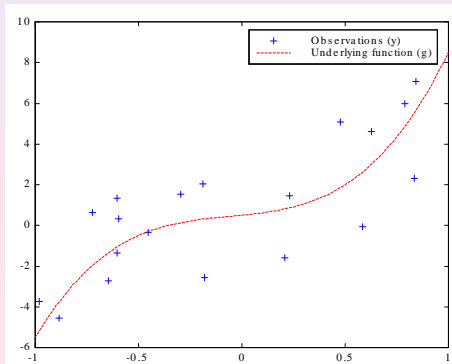$g[(\mathbf{x}(n)]$  True underlying function

$\varepsilon(n)$  i.i.d noise process with zero mean

Example and Task
Idealized regression
Error measures
The real regression
Examples
Model selection

Example
Task
**Data and assumptions**

# Example

Underlying function: $g(x) = 0.5 + x + x^2 + 6x^3$

Noise: $\varepsilon \sim N[0, 2]$

Example and Task
Idealized regression
Error measures
The real regression
Examples
Model selection

Idealized regression
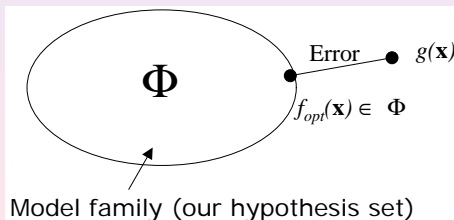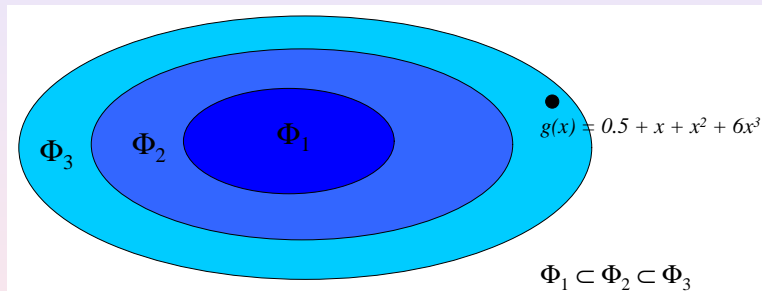
## Idealized regression

Find appropriate model family $\Phi$ and $f(\mathbf{x}) \in \Phi$ with a minimum "distance" (error) to $g(\mathbf{x})$



Model family (our hypothesis set)

Example and Task
**Idealized regression**
Error measures
The real regression
Examples
Model selection

Idealized regression

## Examples of model families



$$\Phi_1 \subset \Phi_2 \subset \Phi_3$$

Linear $\quad \Phi_1 = \{a + bx\}$

Quadratic $\quad \Phi_2 = \{a + bx + cx^2\}$

Cubic $\quad \Phi_3 = \{a + bx + cx^2 + dx^3\}$

Example and Task
**Idealized regression**
Error measures
The real regression
Examples
Model selection

Idealized regression

## How to measure "distance"?

Q: What does the distance between functions $f$ and $g$ mean?

A: The difference between the functions $f$ and $g$.

Q: How do we measure difference (error) between functions?

Example and Task
Idealized regression
**Error measures**
The real regression
Examples
Model selection

The summed squared error (SSE)
Negative log-likelihood
The Bayesian error measure

# The summed squared error (SSE)

$$E = \text{SSE} = \sum_{n=1}^{N} \{f[\mathbf{x}(n), \mathbf{w}] - y(n)\}^2 \tag{1}$$

$\mathbf{w}$ = the parameters of the function $f$.

SSE assumes zero mean i.i.d noise

SSE $\iff$ "Least squares" fit.

Example and Task
Idealized regression
**Error measures**
The real regression
Examples
Model selection

The summed squared error (SSE)
Negative log-likelihood
The Bayesian error measure

## Negative log-likelihood

$$
\begin{align}
\text{Data set } \mathbf{Z} &= \{\mathbf{x}(n), y(n)\}_{n=1,\dots,N} \tag{2} \\
y(n) &= g[(\mathbf{x}(n)] + \varepsilon(n)
\end{align}
$$

$$
E = -\ln L = -\ln \left[ \prod_{n=1}^{N} p[\mathbf{z}(n)|\mathbf{w}] \right] \tag{3}
$$

It is common to assume normally distributed noise $\implies$

$$
p[\mathbf{z}(n)|\mathbf{w}] = p\{f[\mathbf{x}(n), \mathbf{w}] - y(n)\} \sim N[0, \sigma] \tag{4}
$$

This leads to $E \propto \text{SSE}$.

Example and Task
Idealized regression
**Error measures**
The real regression
Examples
Model selection

The summed squared error (SSE)
Negative log-likelihood
The Bayesian error measure

# The Bayesian error measure (1)

- Why maximize the likelihood for the observations given the model parameters?
- Maximize the likelihood for the model parameters given the observations, instead.
- Bayes' theorem tells us how we should do.

Example and Task
Idealized regression
**Error measures**
The real regression
Examples
Model selection

The summed squared error (SSE)
Negative log-likelihood
The Bayesian error measure

# The Bayesian error measure (2)

The probability for the model parameters, given the observations:

$$p(\mathbf{w}|\mathbf{Z}) = \frac{p(\mathbf{Z}|\mathbf{w})p(\mathbf{w})}{p(\mathbf{Z})} = \frac{\mathcal{L}(\mathbf{Z}|\mathbf{w})p(\mathbf{w})}{p(\mathbf{Z})} \tag{5}$$

where $p(\mathbf{w})$ is our "prior" for the model parameters $\mathbf{w}$.
More convenient to minimize the negative likelihood:

$$
\begin{aligned}
E = -\ln p(\mathbf{w}|\mathbf{Z}) &= -\ln \mathcal{L}(\mathbf{Z}|\mathbf{w}) - \ln p(\mathbf{w}) + \ln p(\mathbf{Z}) \\
&\rightarrow = -\ln \mathcal{L}(\mathbf{Z}|\mathbf{w}) - \ln p(\mathbf{w}) \tag{6}
\end{aligned}
$$

since the third term does not depend on the model parameters w.

Example and Task
Idealized regression
Error measures
The real regression
Examples
Model selection

The summed squared error (SSE)
Negative log-likelihood
The Bayesian error measure

## The Bayesian error measure (3)

$$E = -\ln p(\mathbf{w}|\mathbf{Z}) \propto -\ln \mathcal{L}(\mathbf{Z}|\mathbf{w}) - \ln p(\mathbf{w}) \tag{7}$$

Allows including a prior belief, expressed in $p(\mathbf{w})$, about the function $f(\mathbf{x}, \mathbf{w})$.

An example is:

$$p(\mathbf{w}) \propto \exp(-\|\mathbf{w}\|^2/2\sigma_W^2) \tag{8}$$

Example and Task
Idealized regression
**Error measures**
The real regression
Examples
Model selection

The summed squared error (SSE)
Negative log-likelihood
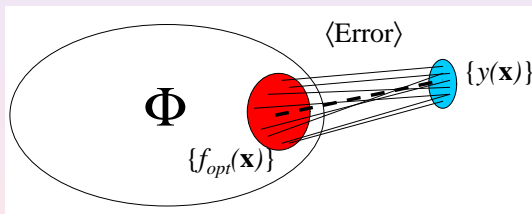The Bayesian error measure

## The Bayesian error measure (4)

- The Bayesian error measure is more general than the ML error.
- The ML error is the special case of the Bayesian error with a uniform prior.
- The Bayesian error is very important to avoid over-fitting.

Example and Task
Idealized regression
Error measures
**The real regression**
Examples
Model selection

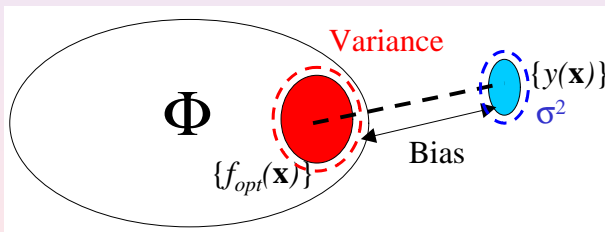The real regression

## The real regression

Find an appropriate model family $\Phi$ and minimize the **expected** distance to $y(\mathbf{x})$ ("generalization error")
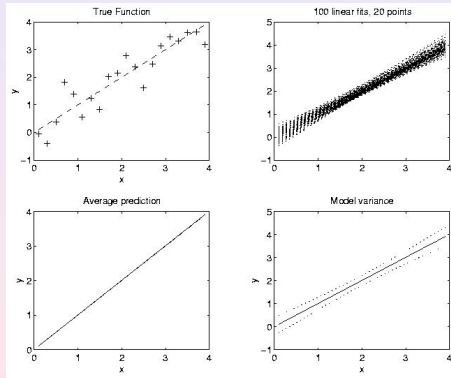


Data is never noise free, and never available in infinite amounts, thus we get variation in data and model. The generalization error is a function of both the training data and the hypothesis selection method.

Example and Task
Idealized regression
Error measures
**The real regression**
Examples
Model selection

The real regression

## Model "bias" & model "variance"

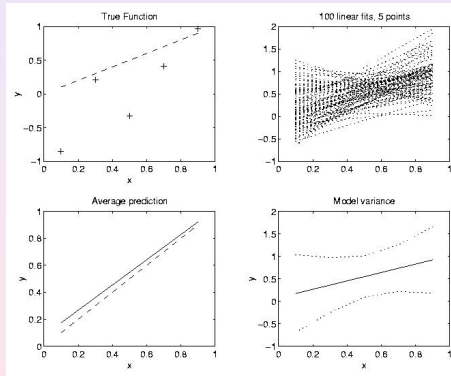$$\langle \mathtt{Error} \rangle = (\mathtt{Bias})^2 + (\mathtt{Variance}) + \sigma_\varepsilon^2 \tag{9}$$
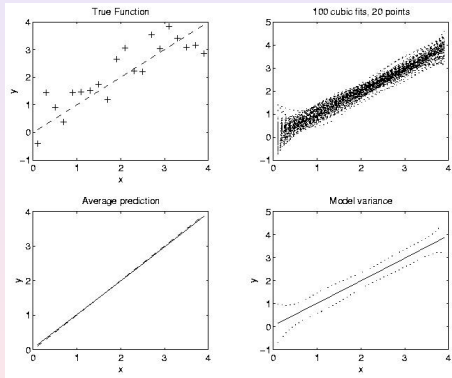
Example and Task
Idealized regression
Error measures
The real regression
Examples
Model selection

Linear function fitted with linear
Linear function fitted with cubic
Quadratic function fitted with linear

# Example (1)



Figure: A linear function $g(x)$ fitted with a linear model $f(x)$, small variance.

Example and Task
Idealized regression
Error measures
The real regression
**Examples**
Model selection

Linear function fitted with linear
Linear function fitted with cubic
Quadratic function fitted with linear

# Example (2)



Figure: A linear function $g(x)$ fitted with a linear model $f(x)$, larger variance.

Example and Task
Idealized regression
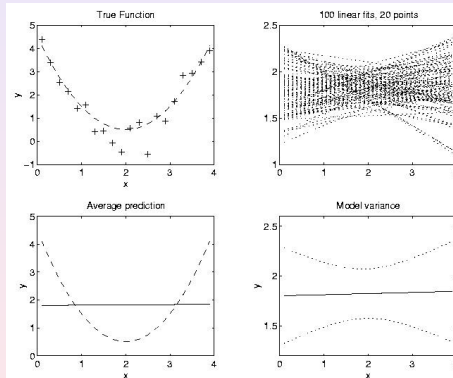Error measures
The real regression
**Examples**
Model selection

Linear function fitted with linear
**Linear function fitted with cubic**
Quadratic function fitted with linear

# Example (1)



Figure: A linear function $g(x)$ fitted with a cubic model $f(x)$, small variance.

Example and Task
Idealized regression
Error measures
The real regression
Examples
Model selection

Linear function fitted with linear
Linear function fitted with cubic
Quadratic function fitted with linear

# Example (2)



Figure: A linear function $g(x)$ fitted with a cubic model $f(x)$, larger variance.

Example and Task
Idealized regression
Error measures
The real regression
**Examples**
Model selection

Linear function fitted with linear
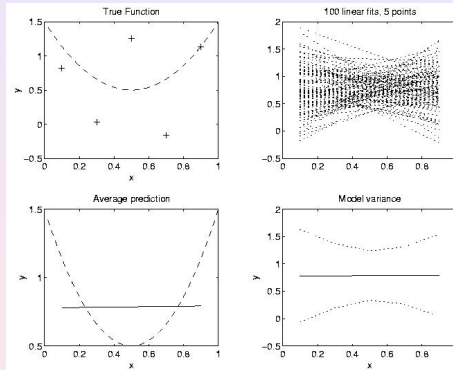Linear function fitted with cubic
**Quadratic function fitted with linear**

# Example (1)



Figure: A quadratic function $g(x)$ fitted with a linear model $f(x)$, small variance.

Example and Task
Idealized regression
Error measures
The real regression
**Examples**
Model selection

Linear function fitted with linear
Linear function fitted with cubic
**Quadratic function fitted with linear**

# Example (2)



Figure: A quadratic function $g(x)$ fitted with a linear model $f(x)$, larger variance.

Example and Task
Idealized regression
Error measures
The real regression
Examples
Model selection

Model selection

# Model selection



Figure: Model with the lowest generalization error is a bias versus variance trade-off.

Example and Task
Idealized regression
Error measures
The real regression
Examples
**Model selection**

Model selection

# Model complexity



Figure: Model with the lowest generalization error is a bias versus variance trade-off.

Example and Task
Idealized regression
Error measures
The real regression
Examples
Model selection

Model selection

## Variable selection

More variables imply larger variance

For linear regression models:

$$\langle E_{\texttt{Test}} \rangle = \langle E_{\texttt{Train}} \rangle + \frac{\sigma_\varepsilon^2 (D+1)}{N} \qquad (10)$$

$\Rightarrow$ A penalty is payed for each input.