

Latent Variable Methods Course

Learning from data

Instructor: Kevin Dunn
kevin.dunn@connectmv.com
<http://connectmv.com>

© Kevin Dunn, ConnectMV, Inc. 2011

Revision: 268:adfd compiled on 15-12-2011

Copyright, sharing, and attribution notice

This work is licensed under the Creative Commons Attribution-ShareAlike 3.0 Unported License. To view a copy of this license, please visit

<http://creativecommons.org/licenses/by-sa/3.0/>



This license allows you:

- ▶ **to share** - to copy, distribute and transmit the work
- ▶ **to adapt** - but you must distribute the new result under the same or similar license to this one
- ▶ **commercialize** - you are allowed to create commercial applications based on this work
- ▶ **attribution** - you must attribute the work as follows:
 - ▶ "Portions of this work are the copyright of ConnectMV", *or*
 - ▶ "This work is the copyright of ConnectMV"

We appreciate:

- ▶ if you let us know about **any errors** in the slides
- ▶ **any suggestions to improve the notes**
- ▶ telling us if you use the slides, especially commercially, so we can inform you of major updates
- ▶ emailing us to ask about different licensing terms

All of the above can be done by writing us at

courses@connectmv.com

If reporting errors/updates, please quote the current revision number: 268:adfd

Overview

Last class: theory, theory, theory.

Today: we cover several concepts and applications

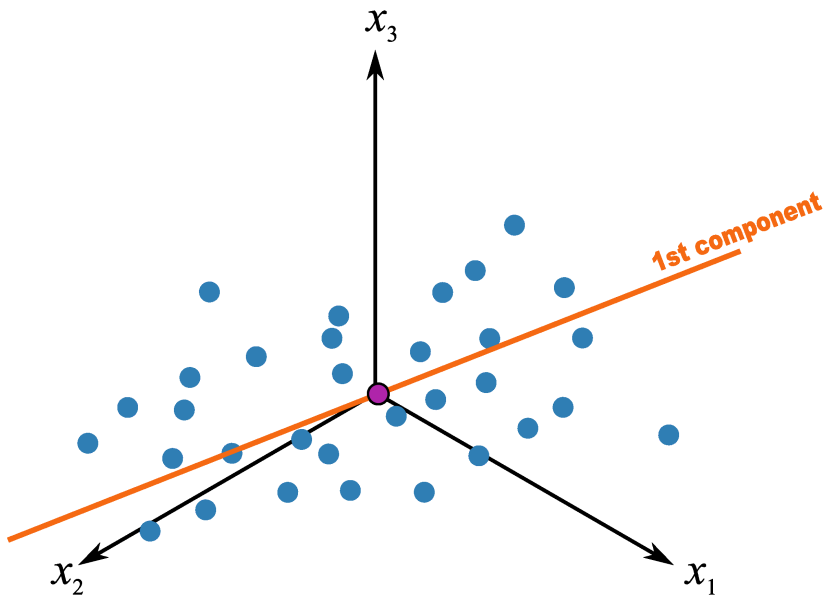
- ▶ how do I use a PCA model after building it?
- ▶ what is overfitting?
- ▶ how many components should I use?
- ▶ various limits used in the PCA model
- ▶ process monitoring with PCA and contribution plots

Last class: Flipping signs

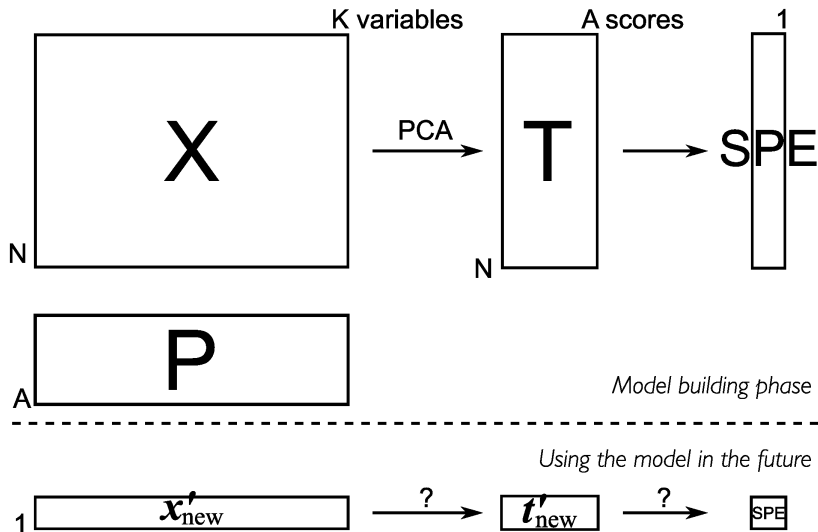
In NIPALS, SVD or eigendecompositions:

- ▶ $\hat{\mathbf{X}}_1 = \mathbf{t}_1 \mathbf{p}'_1 = (-\mathbf{t}_1)(-\mathbf{p}'_1)$
- ▶ Both the scores and loadings may flip sign
- ▶ Depends on the computer, initial guesses, algorithm implementation
- ▶ Not a problem: model interpretation is still consistent
- ▶ Not a problem: model's performance is identical

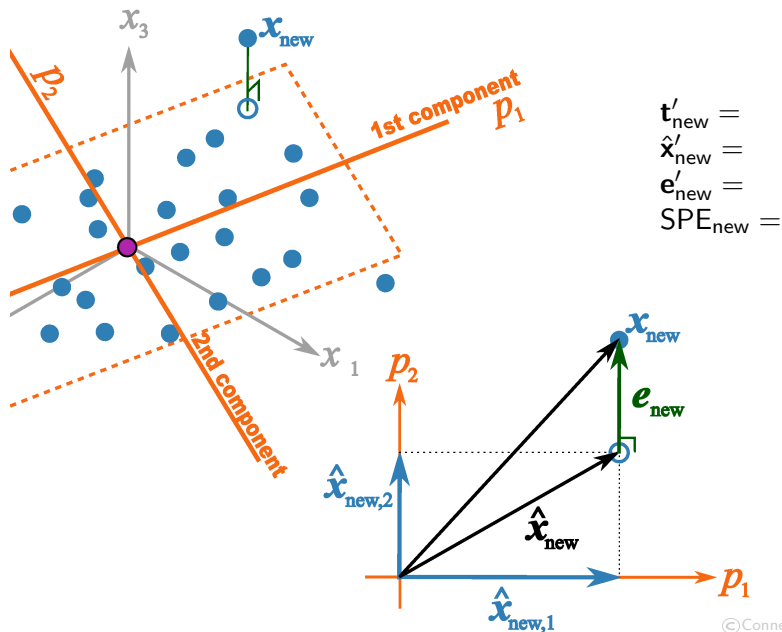
Just be aware when comparing results from different users/software/computers.



Using a PCA model: concept



Using a PCA model: geometric concept



Using existing PCA model on new data

Very simple:

- ▶ Preprocess the raw data: $\mathbf{x}_{\text{new,raw}} \longrightarrow \mathbf{x}_{\text{new}}$
- ▶ Project onto existing model to get scores: $\mathbf{t}'_{\text{new}} = \mathbf{x}'_{\text{new}} \mathbf{P}$
- ▶ Calculated predicted $\hat{\mathbf{x}}'_{\text{new}} = \mathbf{t}'_{\text{new}} \mathbf{P}'$
- ▶ Calculate residuals: $\mathbf{e}'_{\text{new}} = \mathbf{x}'_{\text{new}} - \hat{\mathbf{x}}'_{\text{new}}$
- ▶ Calculate $\text{SPE}_{\text{new}} = \mathbf{e}'_{\text{new}} \mathbf{e}_{\text{new}} = \text{ssq}(\mathbf{e}_{\text{new}})$

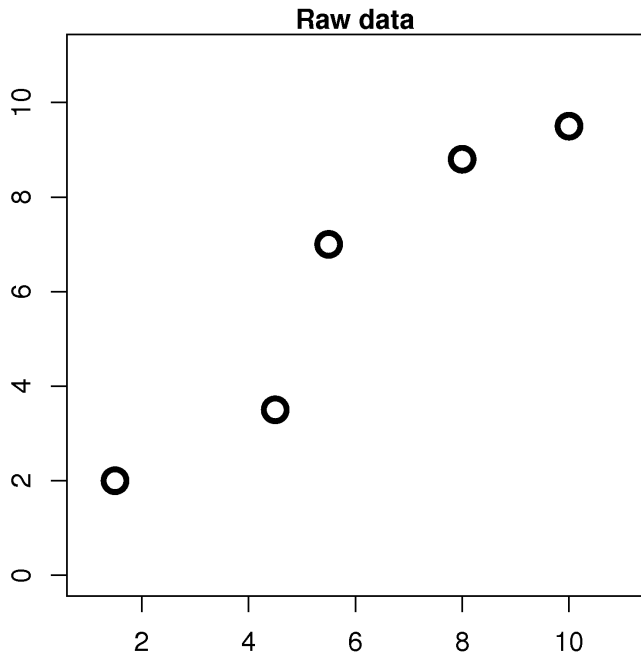
General problem: overfitting

Overfitting

Adding complexity to model when it's not supported **by the data**

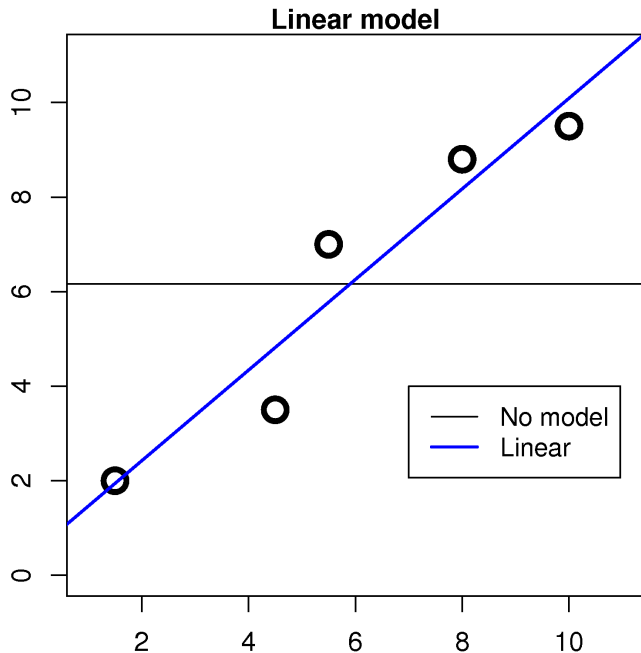
Overfitting is a problem that can occur with *any model*

Overfitting example



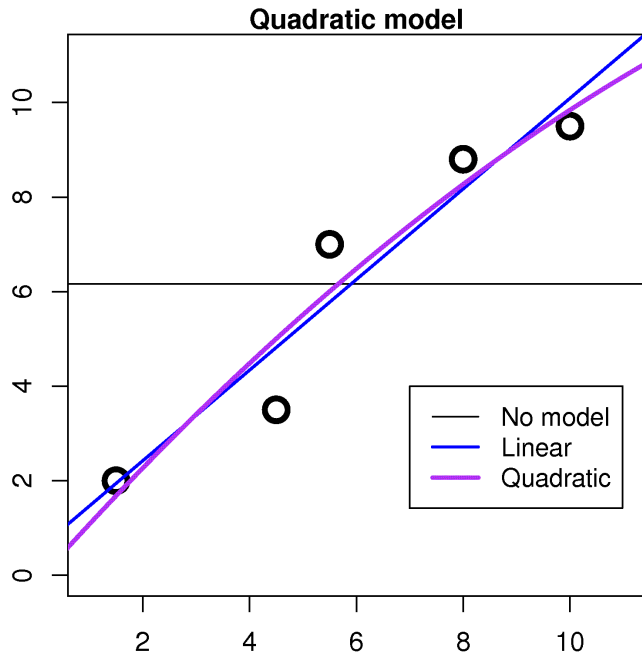
What model
would you fit
initially?

Overfitting example



$$R^2_{\text{linear}} = 0.908$$

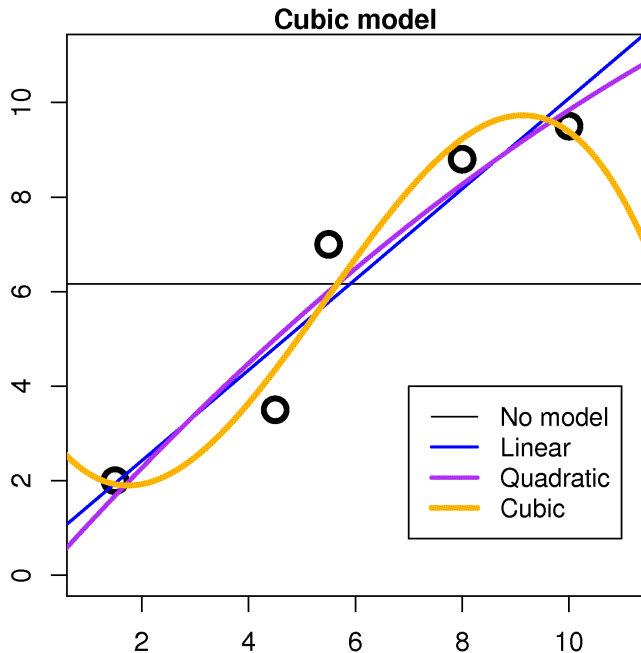
Overfitting example



$$R^2_{\text{linear}} = 0.908$$

$$R^2_{\text{quad}} = 0.913$$

Overfitting example

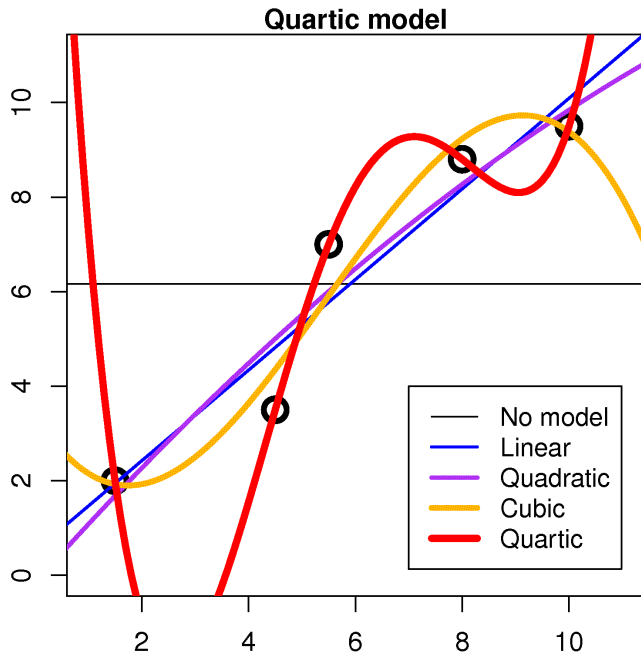


$$R^2_{\text{linear}} = 0.908$$

$$R^2_{\text{quad}} = 0.913$$

$$R^2_{\text{cubic}} = 0.951$$

Overfitting example



$$R^2_{\text{linear}} = 0.908$$

$$R^2_{\text{quad}} = 0.913$$

$$R^2_{\text{cubic}} = 0.951$$

$$R^2_{\text{quartic}} = 1.00$$

Overfitting example

- ▶ **What's the problem here?**

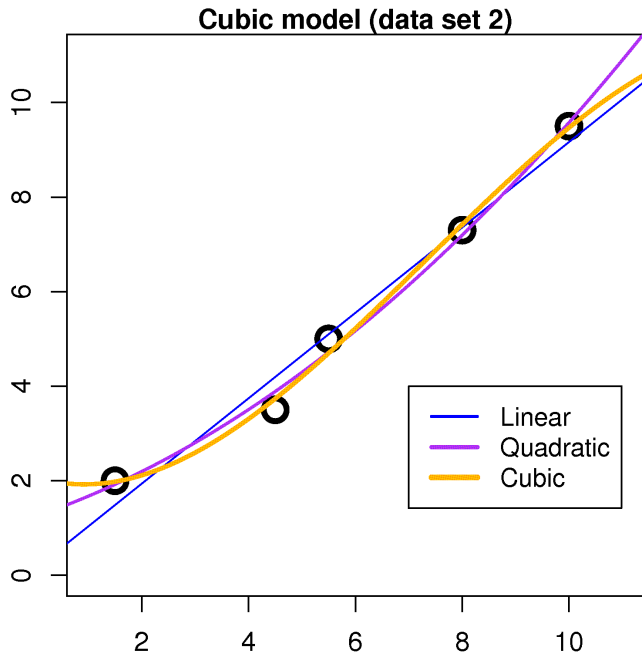
Overfitting issue

Knowing how the model will be used in the future is the key to avoiding overfitting

How many terms would you use in a model if you were:

- ▶ describing the process to a colleague
- ▶ making predictions of y given a new, unknown x
- ▶ making predictions of x for a desired y

Overfitting example 2



Still agree of
your previous
answers?

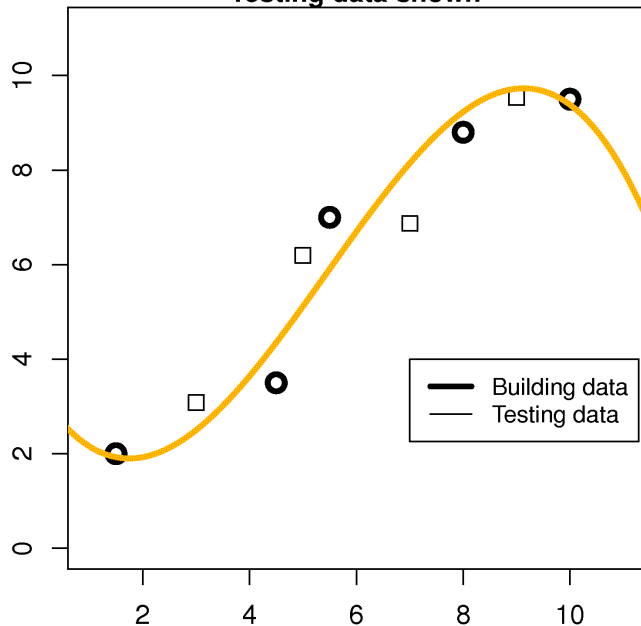
Avoid overfitting: use a testing data set

Principle:

1. Keep a testing data set aside
2. Predict all values from testing data; e.g.
 - ▶ $\hat{y}_i = b_0 + b_1x_i + b_2x_i^2$
3. calculate the residuals
 - ▶ $e_i = y_i - \hat{y}_i$
4. Calculate $\text{ssq}(e_i) = \text{prediction error sum of squares} = \text{PRESS}$

Example: using the testing set

Testing data shown



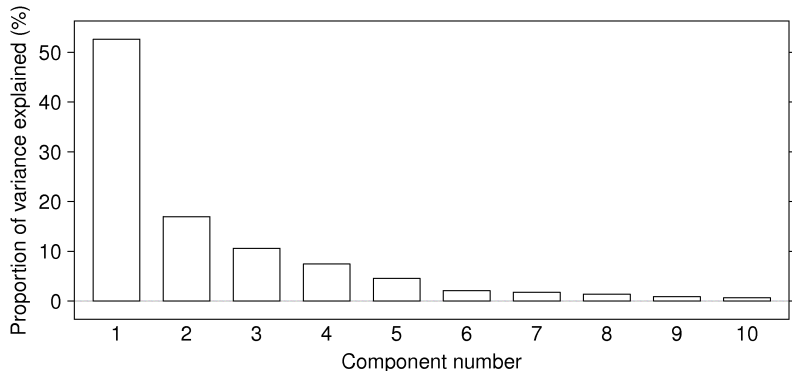
Number of model parameters and PRESS:

1. 21.4 (no model)
2. 1.17 (linear)
3. 1.05 (quadratic)
4. 3.15 (cubic)
5. 29.3 (quartic)

The number of components to use: the problem

- ▶ PCA's objective is to best explain data
- ▶ Over-fitting is when we add more components than are supported by the data in **X**

Pareto plot (scree plot) of R2 per component



How not to do it ...

Fitting to achieve a certain R^2

Should not do this, but it is common to see this approach.

The number of components

Ideal approach

1. Keep a testing data set aside
 2. Fit a component to training data
 3. Project testing data onto model
 4. Calculate the prediction error sum of squares (PRESS) on testing data
 5. Repeat from step 2
-
- ▶ What should happen to PRESS as A increases?
 - ▶ What if we don't have enough data to keep aside?

Cross-validation for PCA

- ▶ A *general* tool; can be applied to any model
- ▶ Cross-validation can help avoid over-fitting
- ▶ $\mathbf{X} = \mathbf{TP}' + \mathbf{E}_A = \mathbf{t}_1\mathbf{p}'_1 + \mathbf{t}_2\mathbf{p}'_2 + \dots + \mathbf{t}_A\mathbf{p}'_A + \mathbf{E}_A$
- ▶ $\mathbf{X} = \mathbf{TP}' + \mathbf{E}_A$

Cross-validation's aim

Find when residuals in \mathbf{E}_A are “*small enough*” so there is no more information left

Cross-validation for PCA

- ▶ $\mathbf{X} = \mathbf{T}\mathbf{P}' + \mathbf{E}_A$
- ▶ $\mathbf{X} = \hat{\mathbf{X}} + \mathbf{E}_A$
- ▶ $\mathcal{V}(\mathbf{X}) = \mathcal{V}(\hat{\mathbf{X}}) + \mathcal{V}(\mathbf{E}_A)$
- ▶ Recall: $R^2 = 1 - \frac{\mathcal{V}(\mathbf{E}_A)}{\mathcal{V}(\mathbf{X})}$
- ▶ Define: $Q^2 = 1 - \frac{\mathcal{V}(\text{predicted } \mathbf{E}_A)}{\mathcal{V}(\mathbf{X})}$
- ▶ $\mathcal{V}(\text{predicted } \mathbf{E}_A) = \text{PRESS} = \text{prediction error sum of squares}$

Cross-validation concept for PCA

Split the rows in \mathbf{X} into G groups.

\mathbf{X}

1	1
2	2
3	3
4	2
5	1
6	3
7	1
8	2
9	3
10	3
11	1
12	2

N

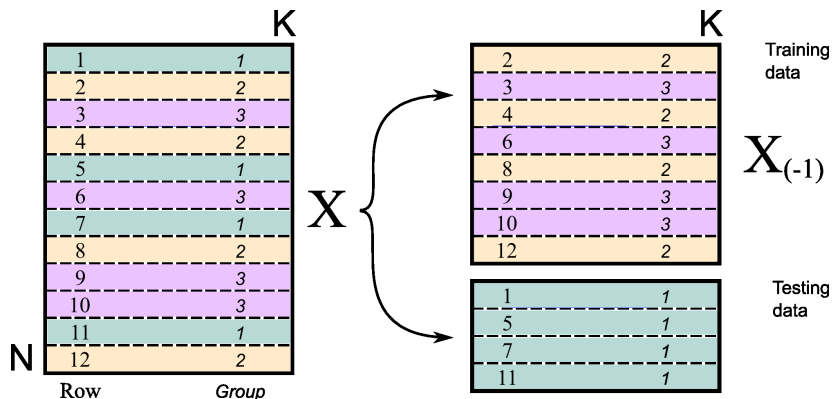
Row Group

- ▶ Typically $G \approx 7$ [ProSensus, Simca-P use $G = 7$]
- ▶ Rows can be randomly grouped, or
- ▶ ordered e.g. 1, 2, 3, 1, 2, 3, ...

$G = 3$ in this illustration

Cross-validation concept for PCA

Fit a PCA model using $\mathbf{X}_{(-1)}$; use $\mathbf{X}_{(1)}$ as testing data

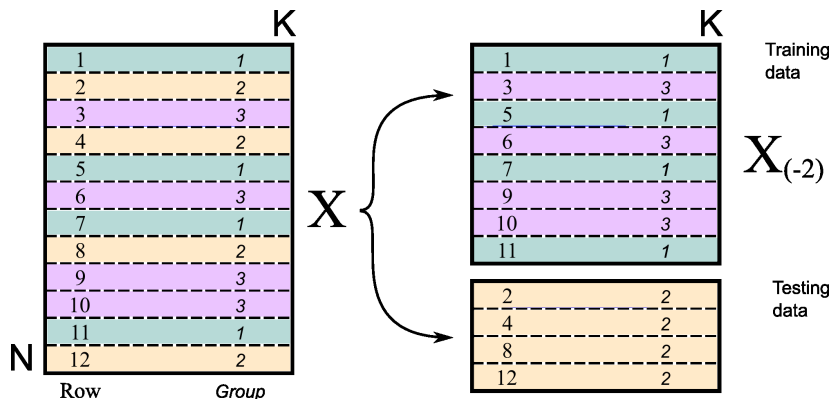


$$\mathbf{E}_{(1)} = \mathbf{X}_{(1)} - \hat{\mathbf{X}}_{(1)}$$

$\mathbf{E}_{(1)}$ = prediction error for testing group 1

Cross-validation concept for PCA

Fit a PCA model using $\mathbf{X}_{(-2)}$; use $\mathbf{X}_{(2)}$ as testing data

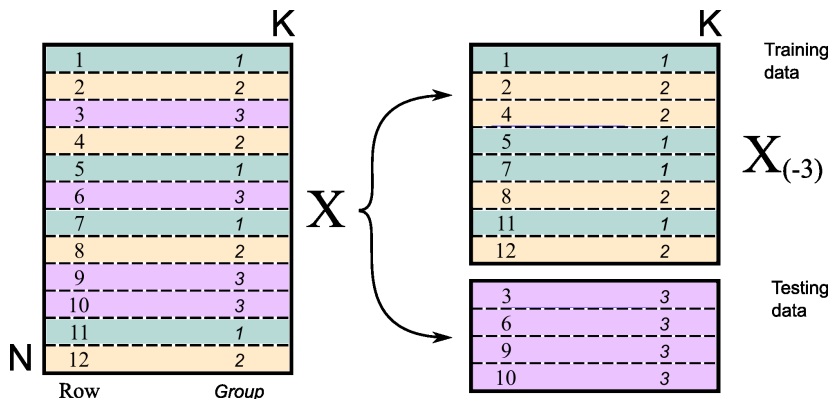


$$\mathbf{E}_{(2)} = \mathbf{X}_{(2)} - \hat{\mathbf{X}}_{(2)}$$

$\mathbf{E}_{(2)}$ = prediction error for testing group 2

Cross-validation concept for PCA

Fit a PCA model using $\mathbf{X}_{(-3)}$; use $\mathbf{X}_{(3)}$ as testing data



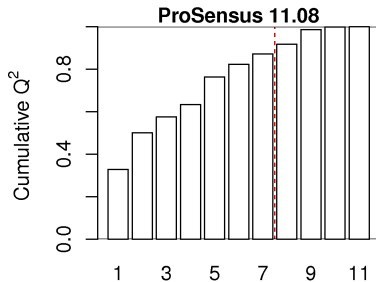
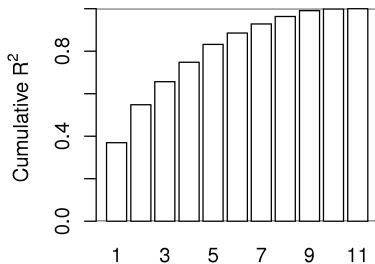
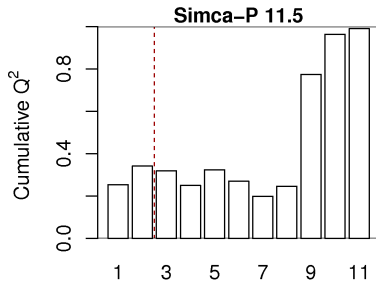
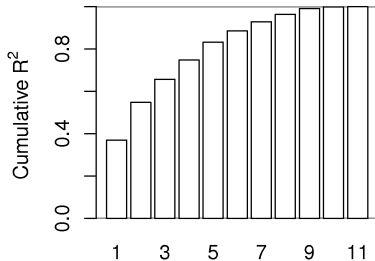
$$\mathbf{E}_{(3)} = \mathbf{X}_{(3)} - \hat{\mathbf{X}}_{(3)}$$

$\mathbf{E}_{(3)}$ = prediction error for testing group 3

Q^2 calculation and interpretation

- ▶ $\text{PRESS} = \text{ssq}(\mathbf{E}_{(1)}) + \text{ssq}(\mathbf{E}_{(2)}) + \dots + \text{ssq}(\mathbf{E}_{(G)})$
- ▶ PRESS = prediction error sum of squares from each prediction group
- ▶ $Q^2 = 1 - \frac{\mathcal{V}(\text{predicted } \mathbf{E}_A)}{\mathcal{V}(\mathbf{X})} = 1 - \frac{\text{PRESS}}{\mathcal{V}(\mathbf{X})}$
- ▶ Q^2 is calculated and interpreted in the same way as R^2
- ▶ Q_k^2 can be calculated for variable $k = 1, 2, \dots, K$
- ▶ You should always find $Q^2 \leq R^2$
- ▶ If $Q^2 \approx R^2$: that component is useful and predictive in the model
- ▶ If Q^2 is “small”: that component is likely fitting noise

Cross-validation: Q^2 can differ in packages



Using default settings in both packages ($G = 7$)

Cross-validation and “autofit”

Each package differs. General idea is to keep the component if

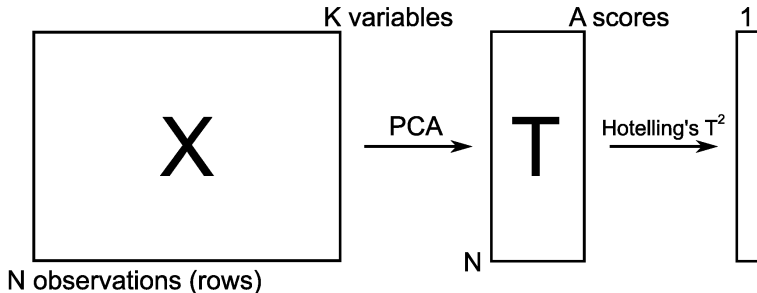
- ▶ Model's Q^2 increases by some percentage (e.g. 1%)
- ▶ Any variable's Q^2 increases by some amount (e.g. 5%)

Summary: number of components to use

- ▶ Cross-validation (autofit) is an “OK” guide: but always test
- ▶ Still an open topic (no single correct method)
 - ▶ Resampling methods are a current research topic
- ▶ Always fit a few extra components in software
- ▶ Do the extra PCs mean anything? Do they help solve your objective? **If so: keep them**

Hotelling's T^2

- ▶ After extracting components from \mathbf{X} we accumulate A score vectors in matrix \mathbf{T}

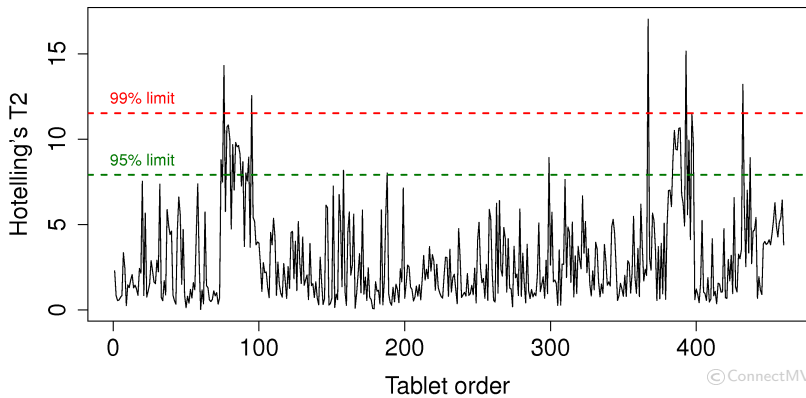


- ▶ T_i^2 is a summary of all A components within row i
- ▶
$$T_i^2 = \sum_{a=1}^{a=A} \left(\frac{t_{i,a}}{s_a} \right)^2$$
- ▶ s_a = standard deviation of score column a

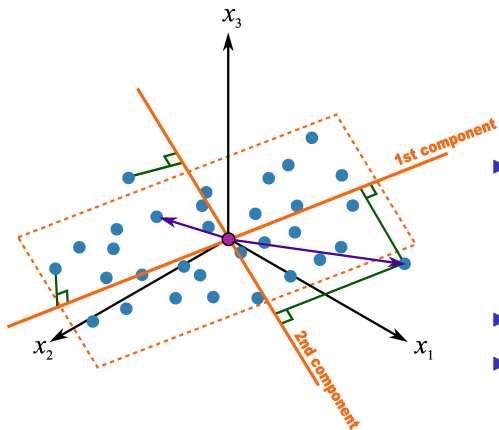
Hotelling's T^2

$$\blacktriangleright T_i^2 = \sum_{a=1}^{a=A} \left(\frac{t_{i,a}}{s_a} \right)^2$$

- ▶ $s_1 > s_2 > \dots$ (from the eigenvalue derivation)
- ▶ $T_i^2 \geq 0$
- ▶ Plotted as a time-series/sequence plot
- ▶ Useful if the row order in dataset has a meaning



Hotelling's T^2

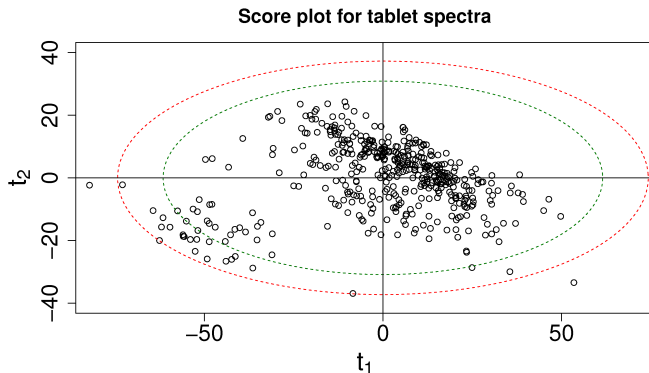


$$T_i^2 = \sum_{a=1}^{a=A} \left(\frac{t_{i,a}}{s_a} \right)^2 \geq 0$$

- Interpretation: directed distance from the center to where the point is projected on the plane
- T^2 has an F -distribution
- Often show the 95% confidence limit value, called $T_{A,\alpha=0.05}^2$

Hotelling's T^2

- ▶ If $A = 2$, equation for 95% limit = $T_{A=2, \alpha=0.05}^2 = \frac{t_1^2}{s_1^2} + \frac{t_2^2}{s_2^2}$
- ▶ An equation for an ellipse
- ▶ s_1 and s_2 are constant for a given model
- ▶ Points on ellipse have a constant distance from model center



Hotelling's T^2

$$T_{A,\alpha}^2 = \frac{(N-1)(N+1)A}{N(N-A)} \cdot F_{\alpha}(A, N-A)$$

- ▶ for A components
- ▶ on N observations
- ▶ for the $100(1-\alpha)\%$ confidence limit

T^2 is also known as:

- ▶ Mahalanobis distance (see [paper on literature website](#))
- ▶ D -statistic

Limits for other PCA model quantities: scores

1. Use a large historical data set (*shown on board*)
2. Use a statistical distribution assumption
 - ▶ $t_a \sim \mathcal{N}(0, s_a)$ (why?)
 - ▶ $100(1 - \alpha)\% \text{limit} = \pm (t_{\alpha/2, df}) s_a$
 - ▶ s_a = standard deviation of score column a
 - ▶ df = degrees of freedom = $N - 1$

Limits for other PCA model quantities: SPE

1. Use a large historical data set (*shown on board*)
2. $SPE_i \sim g\chi^2(h)$ (approximately)
 - ▶ $g = \frac{v}{2m} = \text{premultiplier}$
 - ▶ $h = \frac{2m^2}{v} = \text{degrees of freedom of } \chi^2(h)$
 - ▶ $m = \text{mean}(SPE)$
 - ▶ $v = \text{var}(SPE)$
3. We use the data to estimate g and h
4. See derivation in **Nomikos and MacGregor - paper 34**

Monitoring analogy: your health

- ▶ You have an intuitive (built-in) model for your body
- ▶ When everything is normal: we say “*I’m healthy*” (in control)
- ▶ **Detect a problem:** pain, lack of mobility, hard to breath
- ▶ Something feels wrong (there’s a special cause)
- ▶ **Diagnose the problem:** yourself, search internet, doctor
- ▶ Fix the problem and get back to your usual healthy state

Monitoring analogy: your health

Where did that intuitive model for your body's health come from?

Monitoring analogy: making errors

Assume the doctor is always right and that the baseline hypothesis is: “*you are healthy*”

- ▶ **Type 1 error:** *you detect a problem* (e.g. hard to breathe); doctor says nothing is wrong
 - ▶ You've raised a false alarm
 - ▶ You feel outside your limits,
 - ▶ but the truth is: “*you are healthy*”
 - ▶ **Type 1 error** = raise an alarm when there isn't a problem

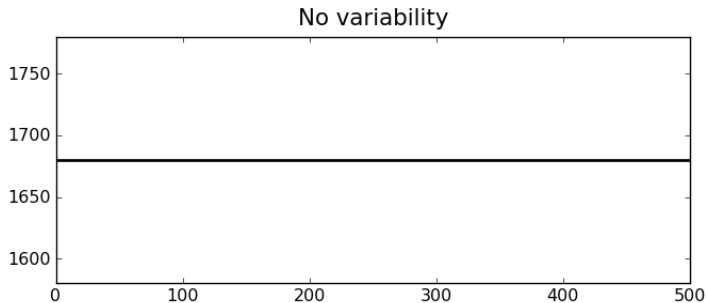
Monitoring analogy: making errors

Assume the doctor is always right and that the baseline hypothesis is: “*you are healthy*”

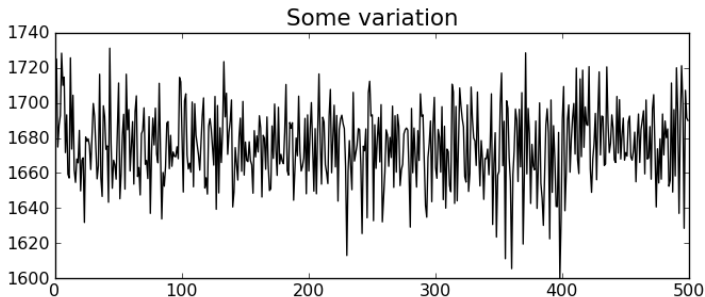
- ▶ **Type 2 error:** *you feel OK*; but go to doctors for physical and they detect a problem
 - ▶ You feel within your limits,
 - ▶ but the truth is: “*you are not healthy*”
 - ▶ **Type 2 error** = don't raise an alarm when there is a problem
- ▶ The grid

Monitoring concept for a process

Our goal: We want process stability



Variability



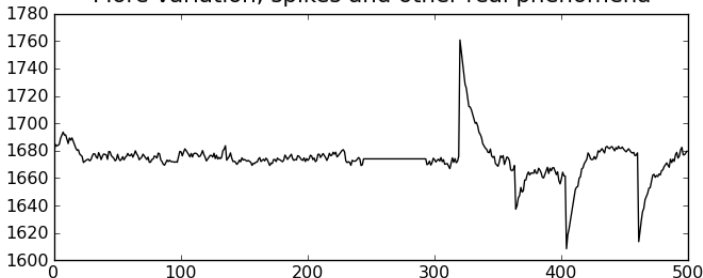
Best case: we have **unaccounted sources** of noise: called **error**

Variability

More realistically:

- ▶ Sensor drift, spikes, noise, recalibration shifts, errors in our sample analysis
- ▶ Operating staff: introduce variability into a process
- ▶ Raw material properties are not constant
- ▶ External conditions change (ambient temperature, humidity)
- ▶ Equipment breaks down, wears out, sensor drift, maintenance shut downs
- ▶ Feedback control introduces variability

More variation, spikes and other real phenomena



Variability in your product

Assertion

Customers expect both uniformity and low cost when they buy your product. Variability defeats both objectives.

Remind yourself of the last time you bought something that didn't work properly

Variability costs you money

The high cost of variability in your final product:

1. Inspection costs:

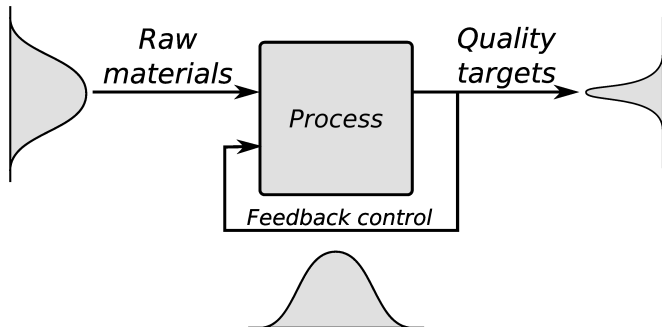
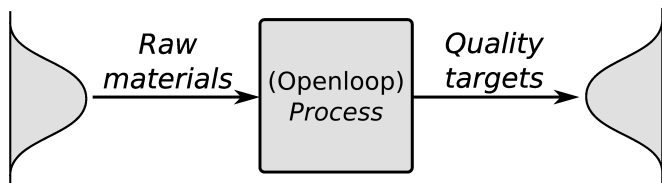
- ▶ high variability: test every product (expensive, inefficient, sometimes destructive)
- ▶ low variability: limited inspection required

2. Off-specification products cost you, and customer, money:

- ▶ reworked
- ▶ disposed
- ▶ sold at a loss

The high cost of variability in your raw materials

- ▶ Flip it around: you receive highly variable raw materials:
 - ▶ That variability lands up in your product, or
 - ▶ you incur additional cost (energy/time/materials) to process it



So what do we want

1. *rapid* problem detection
2. diagnose the problem
3. finally, adjust the process so problems don't occur

Process monitoring is mostly **reactive** and not *proactive*. So it is suited to *incremental* process improvement

Process monitoring: relationship to feedback control

- ▶ “Process monitoring” also called “Statistical Process Control” (SPC)
- ▶ We will avoid this term due to potential confusion:
- ▶ Monitoring is *similar* to (feedback) control:
 - ▶ continually applied
 - ▶ checks for deviations (error)
- ▶ Monitoring is *different* to (feedback) control:
 - ▶ adjustments are **infrequent**
 - ▶ usually **manual**
 - ▶ adjust due to **special causes**
- ▶ Process monitoring: make *permanent* adjustments to reduce variability
- ▶ Feedback control: *temporarily* compensates for the problem

Other types of monitoring you will see

Monitoring is widely used in all industries

- ▶ Managers: monitor geographic regions for hourly sales, downtime, throughput
- ▶ Engineers: monitor large plants, subsections, and unit operations

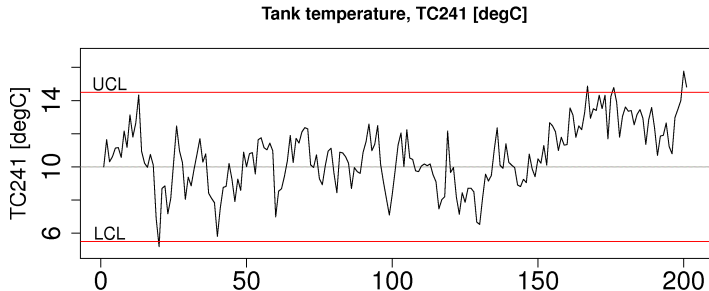
Tools/buzzwords used go by names such as:

- ▶ Dashboards
- ▶ Analytics
- ▶ BI: business intelligence,
- ▶ KPI: key performance indicators

Shewhart chart (recap)

- ▶ Named for *Walter Shewhart* from Bell Telephone and Western Electric, parts manufacturing, 1920's
- ▶ A chart for monitoring variable's *location*, shown with
- ▶ a lower control limit (LCL), usually at $+3\sigma$
- ▶ a upper control limit (UCL), usually at -3σ
- ▶ a target, at the setpoint/desired value

No action taken as long as the variable plotted remains within limits (in-control). Why?



Judging the chart's performance

▶ **Type I error:**

- ▶ value plotted is from common-cause operation, but falls outside limits
- ▶ if values are normally distributed, how many will fall outside?
 - ▶ $\pm 2\sigma$ limits?
 - ▶ $\pm 3\sigma$ limits?
- ▶ *Synonyms*: false alarm, producer's risk

▶ **Type II error:**

- ▶ value plotted is from abnormal operation, but falls inside limits
- ▶ *Synonyms*: false negative, consumer's risk

Adjusting the chart's performance

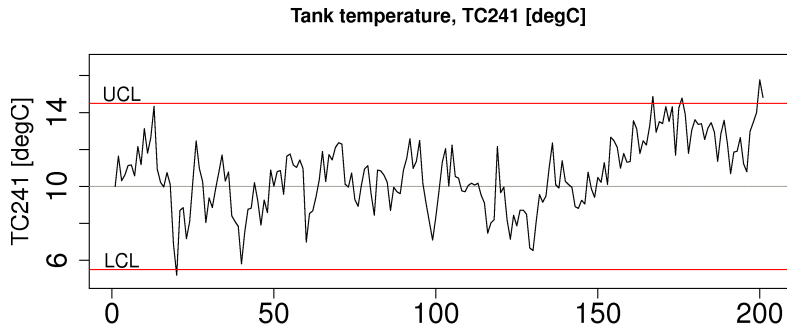
Key point

Control chart limits are not set in stone. Adjust them!

Nothing makes a control chart more useless to operators than frequent false alarms.

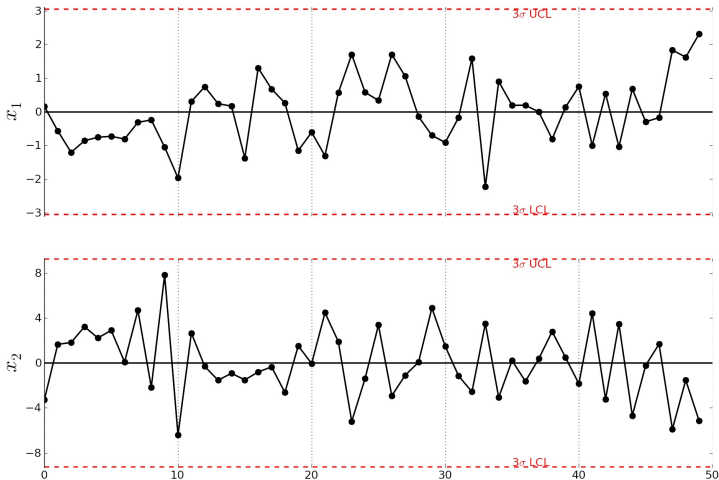
- ▶ **But, you cannot simultaneously have low type I and type II error**

Discussion



1. What action is taken when outside the limits
2. What if data goes missing?

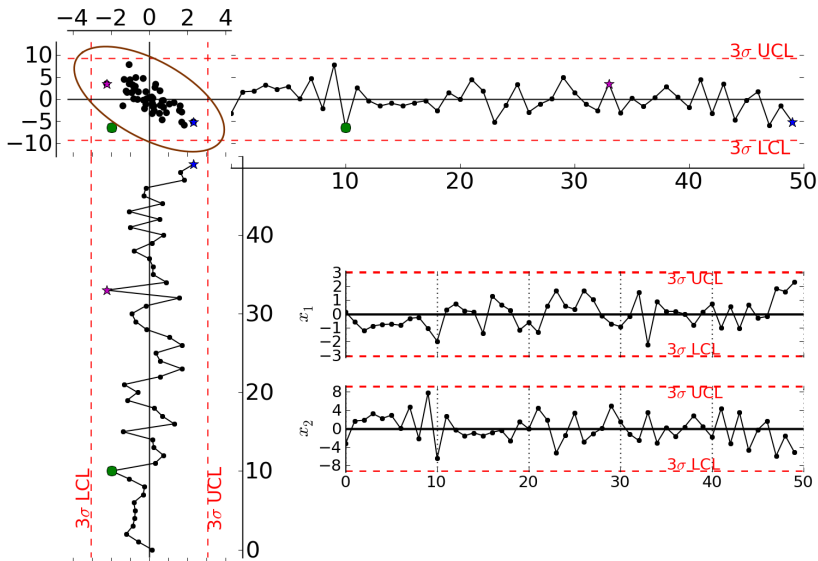
Discussion



3. Monitoring many variables.

- Feasible?
- Is each plot showing something new?

Discussion: multivariate monitoring



Discussion: monitoring only final quality data

Lab measurements have a long time delay:

- ▶ process already shifted by the time lab values detect a problem (continuous)
- ▶ batches have to be placed on hold until lab results return
- ▶ very hard to find cause-and-effect for diagnosis
 - ▶ e.g. low product strength could be caused by multiple reasons

Discussion: monitoring only final quality data

Measurements from real-time systems are:

- ▶ available more frequently (less delay) than lab measurements
- ▶ often are more precise, often with lower error
- ▶ more meaningful to the operating staff
- ▶ contains almost unique “fingerprint” of problem (helps diagnosis)
 - ▶ Now we can figure out what caused low product strength

“Variables” monitored don’t need to be from on-line sensors: could be a calculated value

Process monitoring with PCA: scores

Monitoring with latent variables; use:

- ▶ scores from the model, t_1, t_2, \dots, t_A

Illustration on the board

Process monitoring with PCA: scores

Much better than the raw variables:

- ▶ The scores are orthogonal (independent)
- ▶ Far fewer scores than original variables
- ▶ Calculated even if there are missing data
- ▶ Can be monitored anywhere there is real-time data
- ▶ Available before the lab's final measurement

Process monitoring with PCA: Hotelling's T^2

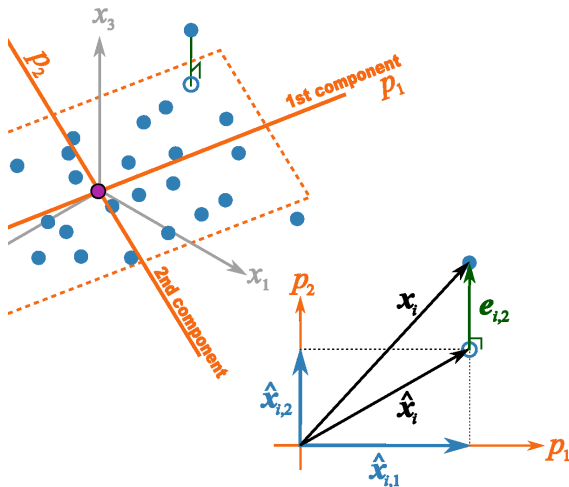
$$\text{Hotelling's } T^2 = \sum_{a=1}^{a=A} \left(\frac{t_a}{s_a} \right)^2$$

- ▶ The distance along the model plane
- ▶ Is a one-side monitoring plot
- ▶ What does a large T^2 value mean?

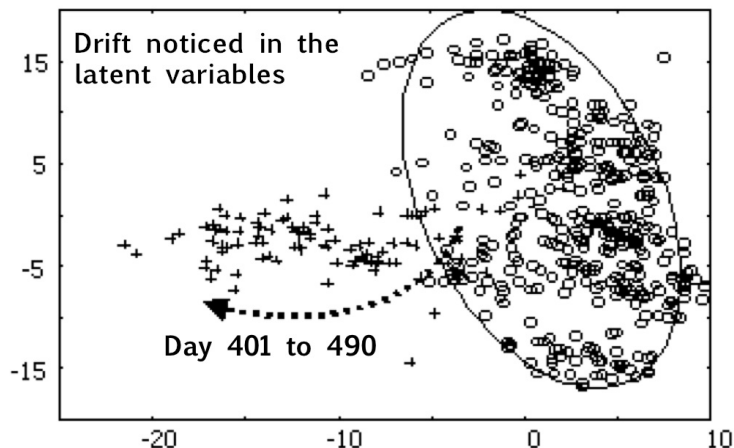
Process monitoring with PCA: SPE

$$\text{SPE}_i = (\mathbf{x}_i - \hat{\mathbf{x}}_i)' (\mathbf{x}_i - \hat{\mathbf{x}}_i) = \mathbf{e}_i' \mathbf{e}_i$$

- ▶ Distance off the model plane
- ▶ Is a one-side monitoring plot
- ▶ What does a large SPE value mean?



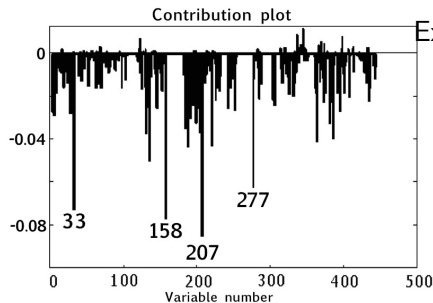
Diagnosing a problem



- Interrogate the latent variables to see what changed

LVM for troubleshooting: contribution plot

- Shows difference between two points in the score plot



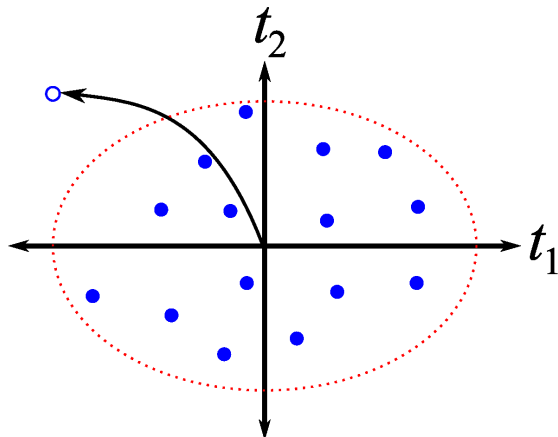
Example:

- **207**: temperature on tray 129 in distillation column 3
 - **158**: a tag from distillation column 3
 - **33** and **277**: related to concentration of feed A
- These variables are related to the problem
 - *Not the cause of the problem*
 - Still have to use your engineering judgement to diagnose
 - But, we've reduced the size of the problem

Contribution plots

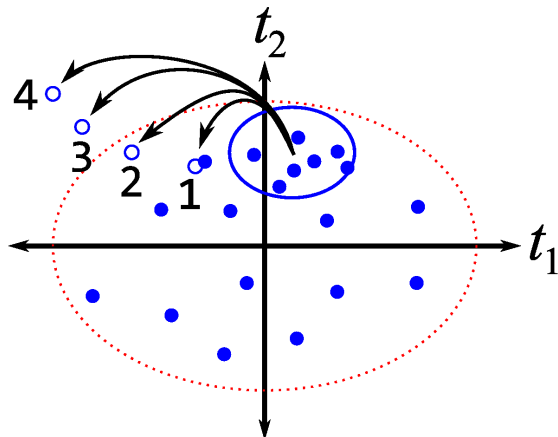
- ▶ Scores: $t_{i,a} = \mathbf{x}_i \mathbf{p}_a$
 - ▶ $[x_{i,1}p_{1,a} \quad x_{i,2}p_{2,a} \quad \dots \quad x_{i,k}p_{k,a} \quad \dots \quad x_{i,K}p_{K,a}]$
 - ▶ *Derivation on the board*
- ▶ T^2 contributions: weighted sum of scores
 - ▶ More details in [Alvarez et al. - paper 21](#)
 - ▶ and [Kourti and MacGregor - paper 81](#)

Contributions in the score space



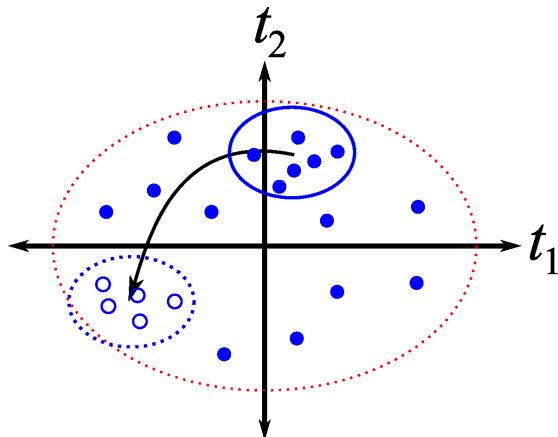
From the model center to a point

Contributions in the score space



Four separate contribution plots to learn why the sequence of deviations occurred

Contributions in the score space



From one group to another group

Contribution plots

- ▶ $SPE = \mathbf{e}_i' \mathbf{e}_i$
 - ▶ where $\mathbf{e}_i' = \mathbf{x}_i' - \hat{\mathbf{x}}_i'$
 - ▶ $[(x_{i,1} - \hat{x}_{i,1}) \quad (x_{i,2} - \hat{x}_{i,2}) \quad \dots \quad (x_{i,K} - \hat{x}_{i,K})]$
- ▶ Joint T^2 and SPE monitoring plots
 - ▶ *Illustrated on the board*
 - ▶ Discussion

Industrial case study: Dofasco

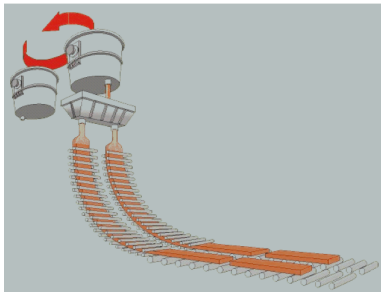
- ▶ ArcelorMittal in Hamilton (formerly called Dofasco) has used multivariate process monitoring tools since 1990's
- ▶ Over 100 applications used daily
- ▶ Most well known is their casting monitoring application, Caster SOS (Stable Operation Supervisor)
- ▶ It is a multivariate monitoring system

Dofasco case study: slabs of steel



All screenshots with permission of Dr. John MacGregor

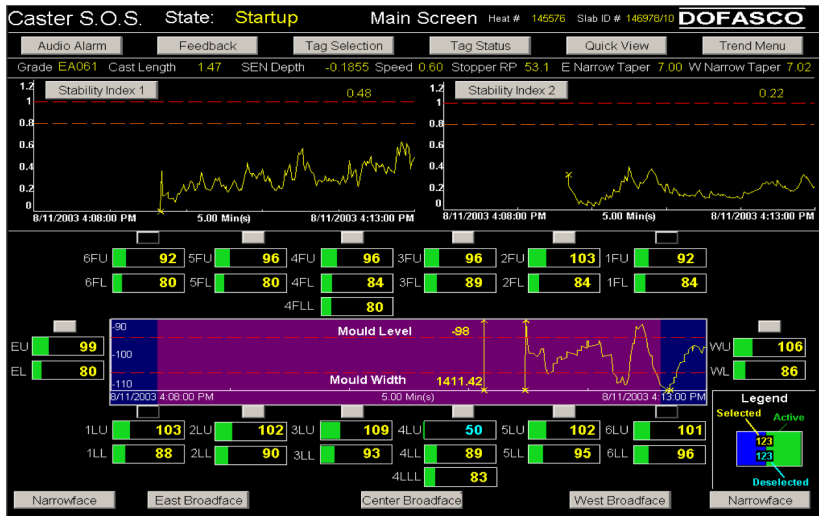
Dofasco case study: casting



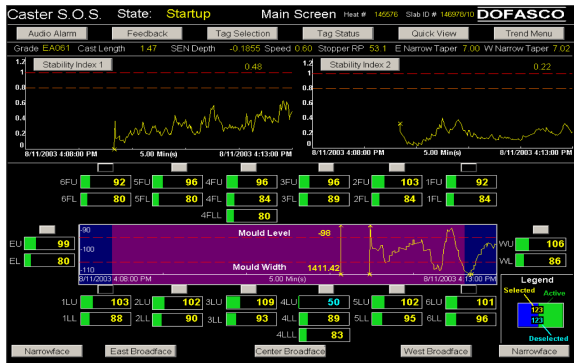
Dofasco case study: breakout



Dofasco case study: monitoring for breakouts

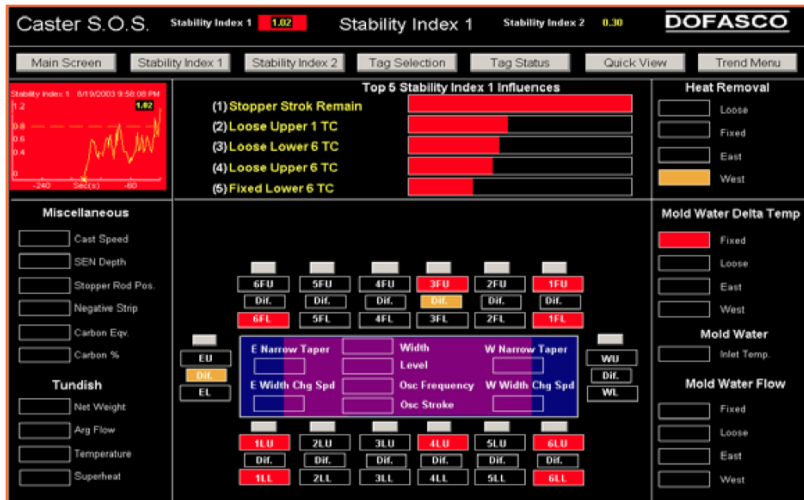


Dofasco case study: monitoring for breakouts



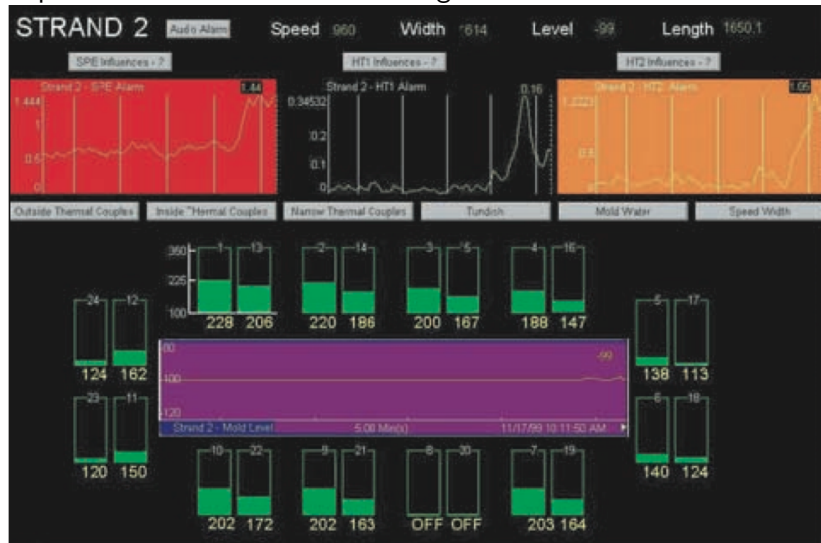
- ▶ Stability Index 1 and 2: one-sided monitoring chart
- ▶ Warning limits and the action limits.
- ▶ A two-sided chart in the middle
- ▶ Lots of other operator-relevant information

Dofasco case study: an alarm



Dofasco case study: previous version

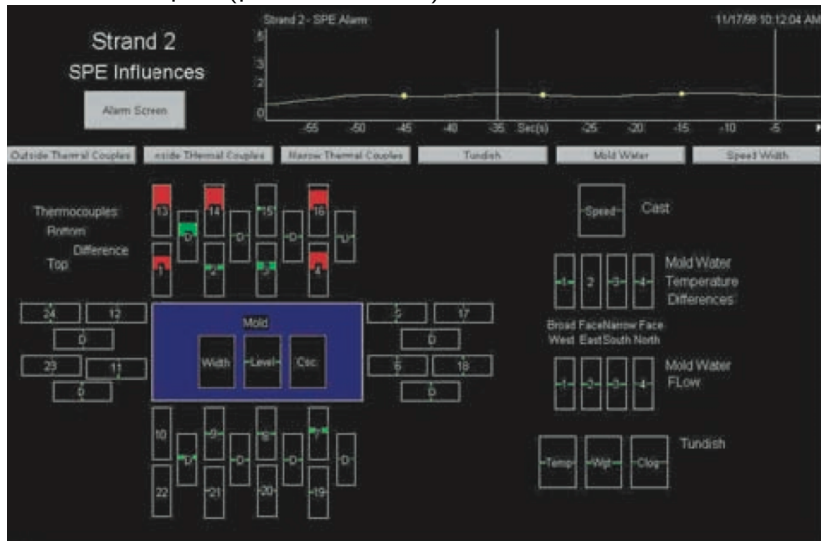
A previous version of the monitoring chart:



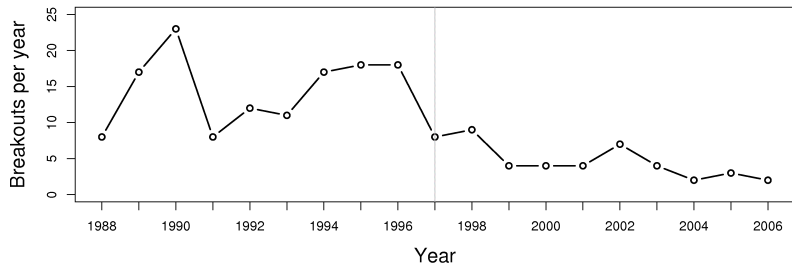
Updated based on operator feedback/requests

Dofasco case study: contribution plots

Contribution plot (previous version):

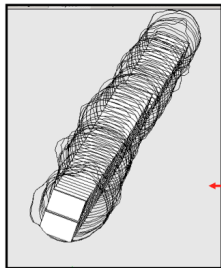


Dofasco case study: economics of monitoring

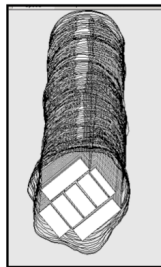


- ▶ Implemented system in 1997; multiple upgrades since then
- ▶ Economic savings: more than \$ 1 million/year
 - ▶ each breakout costs around \$200,000 to \$500,000
 - ▶ process shutdowns and/or equipment damage

Lumber case study

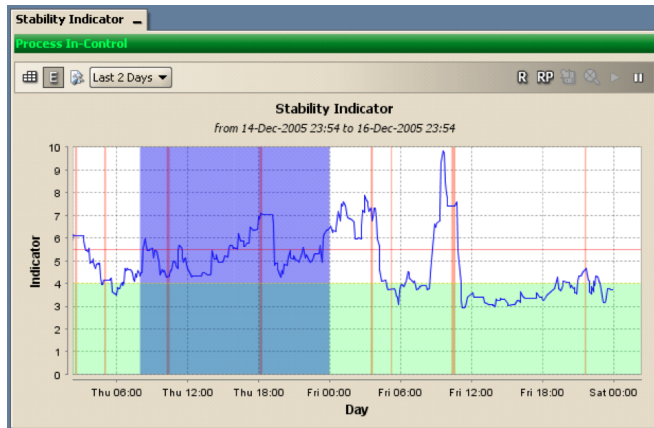


30 to 35 %
50 to 55 %



Show video

Lumber case study



- ▶ Hotelling's T^2 is called "stability indicator" for operators
- ▶ Horizontal red line is the 99% limit
- ▶ Shaded green area is the 0 to 95% limit region

Monitoring isn't just for chemical processes

Any data stream can be monitored

- ▶ Raw material characteristics
- ▶ On-line data from systems (most common multivariate monitoring)
- ▶ Final quality properties
- ▶ End-point detection
- ▶ More generally: **any row in a data matrix**
 - ▶ Credit card/financial fraud monitoring
 - ▶ Human resources

General procedure to build monitoring models I

1. Identify variable(s) to monitor.
2. Retrieve historical data (computer systems, or lab data, or paper records)
3. Import data and just plot it.
 - ▶ Any time trends, outliers, spikes, missing data gaps?
4. Locate regions of stable, common-cause operation.
 - ▶ Remove spikes and outliers
5. Building monitoring model
6. Model includes control limits (UCL, LCL) for scores, SPE and Hotelling's T^2
7. Test your chart on **new, unused** data.
 - ▶ Testing data: should contain both common and special cause operation
8. How does your chart work?
 - ▶ Quantify the type I and II error.

General procedure to build monitoring models II

- ▶ Adjust the limits;
- ▶ Repeat this step, as needed to achieve levels of error
- 9. Run chart on your desktop computer for a couple of days
 - ▶ Confirm unusual events with operators; would they have reacted to it? False alarm?
 - ▶ Refine your limits
- 10. Not an expert system - will not diagnose problems:
 - ▶ use your engineering judgement; look at patterns; knowledge of other process events
- 11. Demonstrate to your colleagues and manager
 - ▶ But go with dollar values
- 12. Installation and operator training will take time
- 13. Listen to your operators
 - ▶ make plots interactive - click on unusual point, it drills-down to give more context

Challenges for real-time monitoring

- ▶ Getting the data out
- ▶ Real-time use of the data (value of data decays exponentially)
- ▶ Training people to use the monitoring system is time consuming
- ▶ Bandwidth/network/storage/computing

Important readings

These papers will help you get to the bottom of process monitoring:

- ▶ MacGregor: **Using on-line process data to improve quality: challenges for statisticians** (paper 75)
- ▶ Kourti and MacGregor: **Process analysis, monitoring and diagnosis, using multivariate projection methods** (paper 31)
- ▶ MacGregor and Kourti: **Statistical process control of multivariate processes** (paper 16)
- ▶ Kresta, MacGregor and Marlin: **Multivariate statistical monitoring of process operating performance** (paper 9)
- ▶ Miller *et al.*: **Contribution plots: a missing link in multivariate quality control** (paper 78)