

**2361-15**

**School on Large Scale Problems in Machine Learning and Workshop on  
Common Concepts in Machine Learning and Statistical Physics**

*20 - 31 August 2012*

**MACHINE LEARNING IN SYSTEMS BIOLOGY: Bioinformatics for Genomic  
Medicine**

Ole WINTHER

*Technical University of Denmark DTU and University of Copenhagen KU  
Denmark*



# Bioinformatics for genomic medicine

Ole Winther

Technical University of Denmark (DTU) & University of Copenhagen (KU)

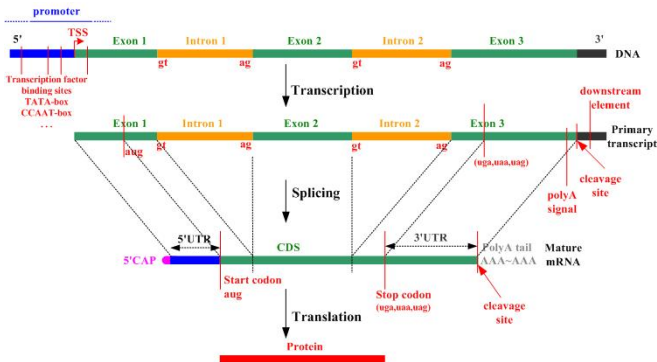
March 28, 2012



- Overall motivation – use genomic data to improve cancer diagnosis and treatment.
- Multidisciplinary collaboration between Bioinformatics Centre, KU and Genomic Medicine and Oncology, Copenhagen University Hospital (Riget).
- Classification and survival analysis for cancer from gene expression data
- Next steps:
  - Include data from more experimental platforms and
  - modeling more aspects of the diagnostic process - therapy selection (personalized medicine)

- **Experimental platforms**
- **Case 1 - Cancer of unknown primary origin (CUP)**
  - Classification of cancer and outlier detection
  - Large gene expression profiling dataset
- **Case 2 - Prognosis in gastric cancer**
  - Random survival forests
  - using gene sets scores as covariates
- Outlook

- **mRNA gene expression profiling**
  - Microarray - 50k genes \$660
- **Exome sequencing** - next generation sequencing
  - [23andme.com/exome](http://23andme.com/exome) 80x coverage of 50M bases \$999
- Additional platforms: Non-coding RNAs and proteomics



- Cancer are classified according to their origin.
- **CUP - a metastasis is located, but not the primary tumor.**
- 2 – 5% of cancer patients get the CUP diagnosis ~ 20 annually at Riget.
- No primary tumor located in two-third of cases.
- Knowing the origin typically determines treatment.
- Cancer is a very heterogeneous disease - improved molecular characterization will lead to identification of more clinical fitting subtypes.

- **Aim: Build classifier for major cancer types**
- Phase one - **data collection and normalization**:
  - Careful curation of 2400+ expression profiles (samples) downloaded from Gene Expression Omnibus (GEO)  
<http://www.ncbi.nlm.nih.gov/geo/>
  - **Training data - 1466 samples**: 1299 primary tumors and 167 normal tissue (various organs)
  - **Test set - 641 tumor samples**: 391 primary tumors and 250 metastases.
  - **57 CUP samples** of which **29** remain unknown after work-up.
- **15 cancer types**: thyroid, lung, stomach, colon/rectum, pancreas, bile duct/gallbladder, liver, kidney, urinary tract, prostate, breast, ovary, endometrium, cervix uteri, testis cancer, a group of malignant melanomas
- **Normalization** with Robust multi-chip average (**RMA**) in **R/Bioconductor package**.

- Phase two - filtering and training of classifier:
  - 47k+ transcript expression values (Affymetrix U133 Plus 2.0)
  - 20k left after variance filtering
- Two-step training of classifier:
  - 1 Univariate test (F) identification of discriminative probes
  - 2 Train classifier on selected probes
- **Optimal p-value cut-off** for F-test ( $= 10^{-180}$ ) found by **nested cross-validation**



- Initial trials identified **linear discriminant analysis (LDA)** as the method with **lowest error rate**
- $c$  is the class label and  $x$  the covariate vector (log fold change in expression values).

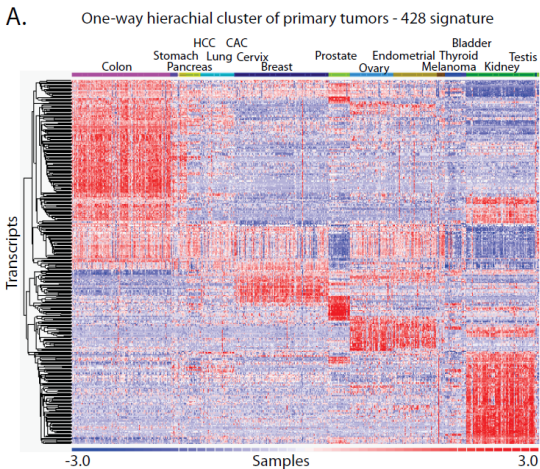
$$p(c|x) = \frac{p(x|c)p(c)}{p(x)}$$

$$p(x|c) = \mathcal{N}(x; \mu_c, \Sigma)$$

$$p(x) = \sum_c p(x|c)p(c)$$

- Test performance 428 probe classifier: **90%** and **83%** for **primary tumors** and **known metastases**, respectively.
- CUP classifier (merge training and test set, 641 probes) had a **LOOCV** accuracy of **92%** and **87%**, respectively.

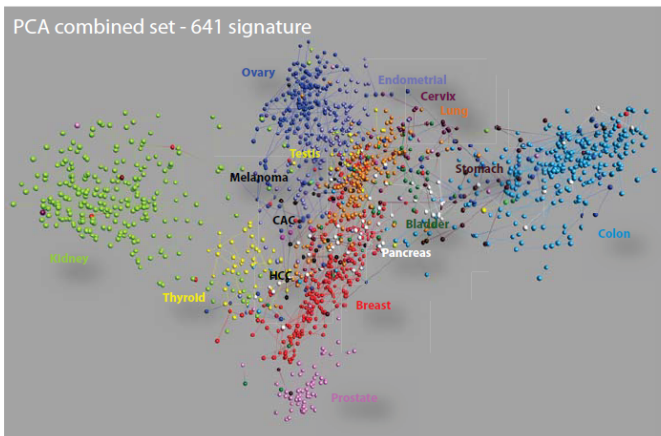
## Cancer of unknown primary origin (CUP) classification





## Cancer of unknown primary origin (CUP) classification

B.





## Classification of CUP patients

ID	Sex/age	Biopsy site	Histology	Path Diag.	Stand of Ref	LDA Pred	Outlier score
14.	F/56	LN neck	PDC	Lung	Lung (CD)	Lung	975
17.	F/57	LN neck	Adenoc.	Lower GI	Colon (RD)	Colon	746
22.	M/55	LN neck	Adenoc.	CUP	Stomach (RD)	Normal	934
23.	M/39	LN retro	PDC	CUP	Kidney (RD)	Kidney	1085
28.	F/58	Peritoneum	PDA	Ovary	Ovary (RD)	Ovary	810
31.	F/40	LN neck	PDA	CUP	Lung (CD)	Stomach	985
34	M/74	Skin	PDA	Lung	Lung (RD)	Lung	898
39.	M/71	Liver	Adenoc.	CUP	Pancreas (CD)	Pancreas	1097
40.	F/44	Liver	Adenoc.	Colon	Colon (RD)	Colon	729
44.	F/43	Kidney	Carc.	CUP	Bladder (RD)	Bladder	1286
49.	M/60	LN neck	PDA	Kidney	Kidney (RD)	Kidney	1223
51.	F/42	LN pelvis	SCC	CUP-SCC	Cervical (RD)	Cervix	828
52.	M/53	Liver	PDA	CUP	CCC (RD)	CCC	923
53.	M/70	Liver	Adenoc.	Lung	Lung (RD)	Lung	1047
57.	M/67	Liver	Adenoc.	CCC	CCC (RD)	HCC	965
66.	F/68	Liver	PDA	CUP	CCC (RD)	Cervix	1100
70.	M/38	Peritoneum	Adenoc.	Stomach	Stomach (CD)	Colon	842
74.	M/62	Leg	Carc.	Adnex tumor	Adnex tumor (RD)	Normal	1010
76.	M/64	Liver	Adenoc.	Lower GI	Small intestine (RD)	Colon	912
77.	M/59	LN axilla	PDC	CUP	Lung (CD)	Breast	978
86.	F/61	LN axilla	Adenoc.	CUP	Lung (RD)	Stomach	1108
88.	F/36	Peritoneum	Adenoc.	Ovary	Ovary (RD)	Cervix	1033
89.	F/57	Liver	PDA	CCC	CCC (RD)	CCC	916
90.	F/71	Peritoneum	Adenoc	Ovary	Ovary (RD)	Ovary	781
92.	M/62	Liver	Malignant tumor	Angiosarcoma	Angiosarcoma (RD)	Normal	1097
95.	M/45	Peritoneum	PDC	DSRCT	DSRCT (RD)	Breast	1098
71+72	M/61	Bone + Kidney	PDC	Kidney	Kidney (RD)	Kidney	1096
75+87	F/43	Liver	PDA	CCC	CCC (RD)	CCC	1277
							925
							1030



## Classification of CUP patients

ID	Sex/age	Biopsy site	Histology	Path Diag.	Stand of Ref	LDA Pred	Outlier score
11.	F/58	LN neck	PDA	CUP	CUP (SD)	Ovary	756
13.	F/72	Peritoneum	PDA	CUP	CUP (NSD)	Pancreas	1193
21.	M/63	LN neck	PDC	CUP	CUP (NSD)	Breast	1108
26.	F/67	Skin	PDA	CUP	CUP (NSD)	Breast	971
32.	M/53	LN neck	PDSCC	CUP-SCC	CUP (NSD)	Normal	926
33.	M/58	Skin	PDA	CUP	CUP (NSD)	Colon	1098
41.	M/74	Liver	PDA	Pancreas	CUP (NSD)	Stomach	1040
42.	M/56	Liver	Adenoc.	CUP	CUP (NSD)	Pancreas	994
43.	F/50	LN retro	PDA	CUP	CUP (NSD)	Stomach	797
45.	M/44	Liver	PDC	CUP	CUP (NSD)	Colon	1245
46.	F/76	Liver	Adenoc.	CUP	CUP (NSD)	Normal	1027
47.	F/59	Liver	Adenoc.	CUP	CUP (SD)	CCC	932
48.	F/59	LN neck	PDC	CUP	CUP (NSD)	Ovary	1032
54.	F/67	Liver	Adenoc.	CUP	CUP (NSD)	Normal	1068
55.	F/55	Liver	Adenoc.	CUP	CUP (NSD)	Normal	962
58.	F/67	Liver	PDC	CUP	CUP (SD)	CCC	995
61.	M/72	Liver	Carc.	HCC	CUP (NSD)	CCC	1102
64.	F/65	LN inguini	PDA	CUP	CUP (SD)	Lung	1168
65.	M/62	LN neck	PDSCC	CUP-SCC	CUP (NSD)	Breast	929
73.	M/43	LN retro	PDC	CUP	CUP (NSD)	Normal	1020
78.	F/59	Lung	Adenoc.	Lower GI	CUP (NSD)	Lung	1062
80.	F/58	Liver	Adenoc.	CUP	CUP (SD)	CCC	1111
81.	F/71	Liver	PDA	CUP	CUP (NSD)	Breast	1212
82.	F/56	Bone	Adenoc.	CUP	CUP (NSD)	CCC	1209
83.	F/59	Liver	PDA	CUP	CUP (SD)	CCC	1061
91.	F/65	LN axilla	Adenoc.	CUP	CUP (SD)	Lung	939
93.	M/58	Bone	PDSCC	CUP-SCC	CUP (NSD)	Breast	940
94.	F/55	Liver	PDA	CUP	CUP (NSD)	Normal	984
50. + 68	M/41	Adr gl	PDC	CUP	CUP (NSD)	Stomach	978
						Pancreas	1079

- CUP classifier not sex-specific, but some cancers are: ovary, cervical and prostate.
- **Unlikely events occur** - like men classified as breast - **may reflect real biology and limited validity of classification categories.**
- Renormalizing class priors  $p(c)$ :

$$p(c|x) = \frac{p(x|c)p(c)}{p(x)}$$

$$p_{\text{rm}}(c|x) = \frac{p(x|c)p_{\text{rm}}(c)}{p_{\text{rm}}(x)} = p(c|x) \frac{p(x)}{p_{\text{rm}}(x)} \frac{p_{\text{rm}}(c)}{p(c)}$$

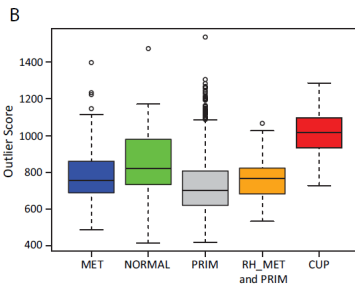
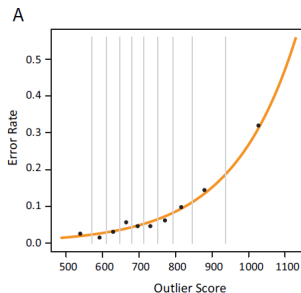
- $p_{\text{rm}}(c) = 0$  for cancers not appearing (like ovary in men) and renormalize others appropriately.

- $p(x)$  measure how probable  $x$  is according to the model - **outlier score**

$$\text{OS} = -\log p(x)$$

- Use **quadratic discriminant analysis (QDA)**

$$p(x) = \sum_c p(x|c)p(c) \quad \text{with} \quad p(x|c) = \mathcal{N}(x; \mu_c, \Sigma_c)$$







# Probabilistic PCA

- Probabilistic PCA (Tipping and Bishop, 1999):

$$x = Wz + \epsilon$$

$$z \sim \mathcal{N}(z; 0, I)$$

$$\epsilon \sim \mathcal{N}(\epsilon; 0, \sigma^2 I)$$

- Marginalizing  $z$  and  $\epsilon$

$$p(x|W, \sigma^2) = \mathcal{N}(x; 0, WW^T + \sigma^2 I)$$

- **Structured covariance** model.

- Log likelihood for  $\mathbf{W}$  and  $\sigma^2$ :

$$\begin{aligned}\log L(\theta; \mathbf{X}) &= \sum_n \log p(x_n | \mathbf{W}, \sigma^2) \\ &= -\frac{N}{2} \left\{ \log \det 2\pi \Sigma + \text{Tr} \left[ \Sigma^{-1} \mathbf{S} \right] \right\}\end{aligned}$$

- Model covariance:  $\Sigma = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I}$
- Empirical covariance:  $\mathbf{S} = \frac{1}{n} \mathbf{X}^T \mathbf{X}$
- Spectral decomposition:  $\mathbf{S} = \mathbf{U} \Lambda \mathbf{U}^T$ ,  $\Lambda_{ii} \geq \Lambda_{jj}$  for  $i < j$ .
- Maximum likelihood solution:**  $i = 1, \dots, m$

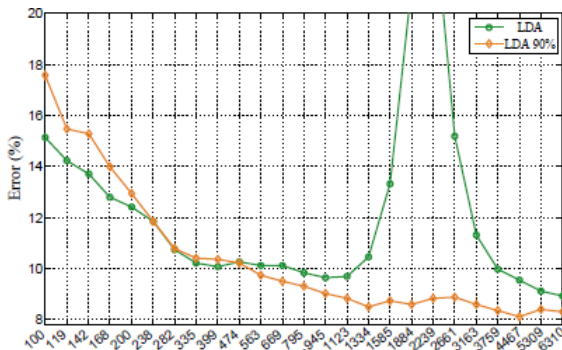
$$w_{i,\text{ml}} = u_i \sqrt{\Lambda_{ii} - \sigma_{\text{ml}}^2} R \quad \text{and} \quad \sigma_{\text{ml}}^2 = \frac{1}{p-m} \sum_{i=m+1}^p \Lambda_{ii}$$

- $R$  arbitrary rotation

- Error versus  $p$  for fixed “variance explained”

$$1 - \frac{1}{\text{Tr} \Lambda} \sum_{i=m+1}^p \Lambda_{ii}$$

- Covariates sorted according to variance.



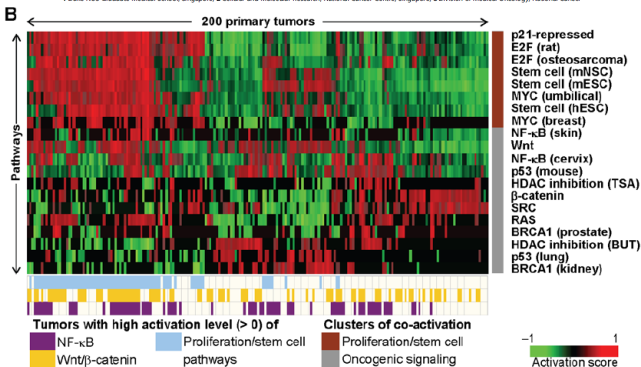
OPEN ACCESS [Freely available online](#)

PLOS GENETICS

## Oncogenic Pathway Combinations Predict Clinical Prognosis in Gastric Cancer

Chia Huey Ooi<sup>1</sup>, Tatiana Ivanova<sup>2</sup>, Jeanie Wu<sup>2</sup>, Minghui Lee<sup>2</sup>, Iain Beehuat Tan<sup>3</sup>, Jiong Tao<sup>2,4</sup>, Lindsay Ward<sup>5</sup>, Jun Hao Koo<sup>2</sup>, Veena Gopalakrishnan<sup>2</sup>, Yansong Zhu<sup>2</sup>, Lai Ling Cheng<sup>6</sup>, Julian Lee<sup>2</sup>, Sun Young Rha<sup>7</sup>, Hyun Cheol Chung<sup>7</sup>, Kumaresan Ganesan<sup>2</sup>, Jimmy So<sup>8</sup>, Khee Chee Soo<sup>9</sup>, Dennis Lim<sup>10</sup>, Weng Hoong Chan<sup>10</sup>, Wai Keong Wong<sup>10</sup>, David Bowtell<sup>11</sup>, Khay Guan Yeoh<sup>12</sup>, Heike Grabsch<sup>5</sup>, Alex Boussioutas<sup>11,13</sup>, Patrick Tan<sup>1,2,14,15\*</sup>

<sup>1</sup> Duke-NUS Graduate Medical School, Singapore, <sup>2</sup> Cellular and Molecular Research, National Cancer Centre, Singapore, <sup>3</sup> Division of Medical Oncology, National Cancer



- This study inspired us to consider survival analysis and gene expression data.
- Use **gene set activation scores** instead of gene expression.
- **Curated gene sets from Msig database**  
[www.broadinstitute.org/gsea/msigdb/](http://www.broadinstitute.org/gsea/msigdb/) and scoring according to GAGE (Luo et al, 2009):

$$t = \frac{m - M}{\sqrt{s^2/n + S^2/n}}$$

- Instead of using a preselected list of gene sets
- we use unbiased search with **random survival forest** among all gene sets.

- Random survival forest identifies gene sets which makes biological sense
- according to literature and may be validated by clinicians:

[1] "OHM\_EMBRYONIC\_CARCINOMA\_UP\_DN\_ud"

[2] "ST\_FAS\_SIGNALING\_PATHWAY\_b"

[3] "TONKS\_TARGETS\_OF\_RUNX1\_RUNX1T1\_FUSION\_SUSTAINED\_IN\_MONOCYTE\_UP\_DN\_ud"

[4] "WEIGEL\_OXIDATIVE\_STRESS\_BY\_HNE\_AND\_H2O2\_u"

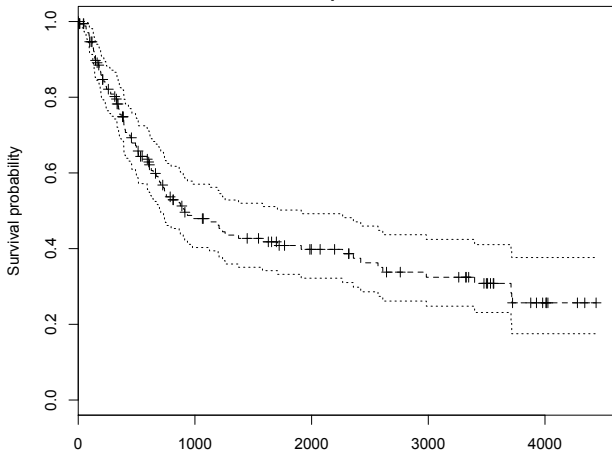
[5] "BIOCARTA\_IL1R\_PATHWAY\_b"

[6] "SNIJDERS\_AMPLIFIED\_IN\_HEAD\_AND\_NECK\_TUMORS\_u"

[7] "GARGALOVIC\_RESPONSE\_TO\_OXIDIZED\_PHOSPHOLIPIDS\_LIGHTYELLOW\_UP\_DN\_ud"

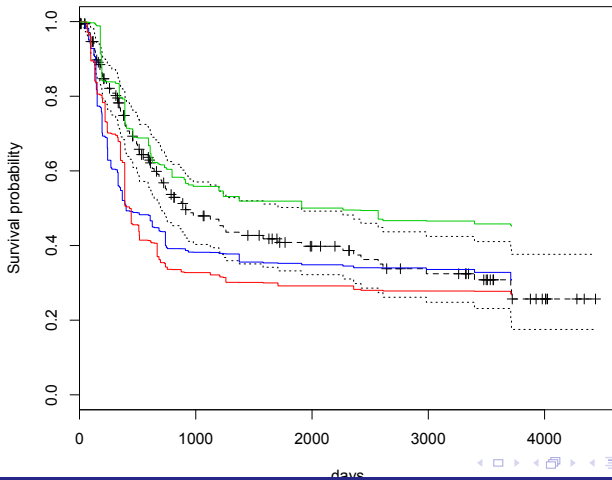


**Kaplan-Meier estimate for overall survival  
(with error curves)  
177 patients**





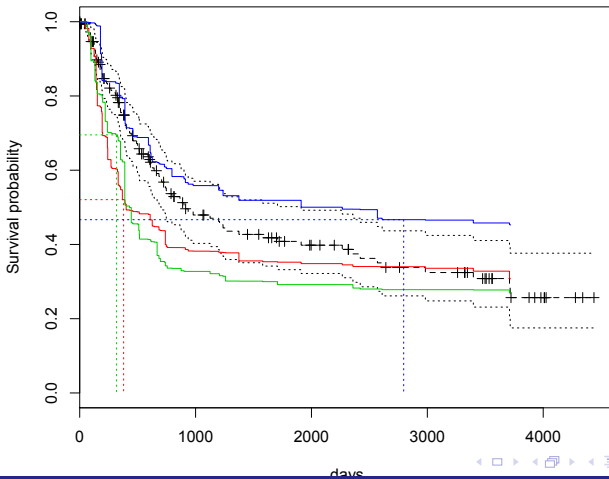
### Predicted individual survival curves for 3 new patients





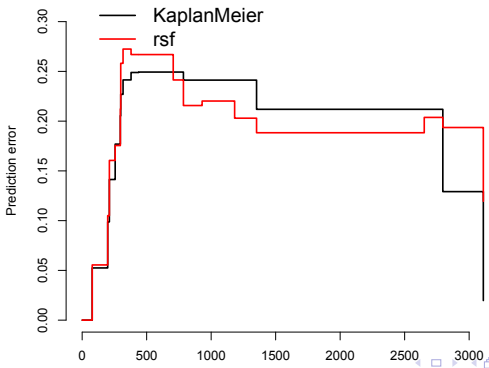


### Predicted individual survival curves for 3 new patients

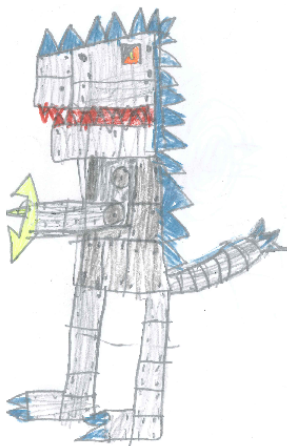




- Hold-out (19 of 198) Brier scores
- Kaplan Meier 0.201
- RSF 0.195



- Proof of concept of close collaboration with clinicians
- Next steps towards individualized treatment
  - better models
  - richer genomic data
  - predict response to treatment
- International Genomics Consortium (ICG) and The Cancer Genome Atlas (TCGA) provide unprecedented amount of genomic data, but close collaboration with clinicians necessary to get detailed clinical data.





## Bioinformatics

- Bogumil Kaczkowski
- Ricardo Henao
- Tomas Martin-Bertelsen
- Anders Krogh

## Genomic Medicine, Riget

- Finn Cilius Nielsen
- Lennart Friis-Hansen
- Rehannah Borup

## Oncology, Riget

- Gedske Daugaard
- Anne Kirstine Moller