

My first scientific paper
Week 5
Highlight the principles

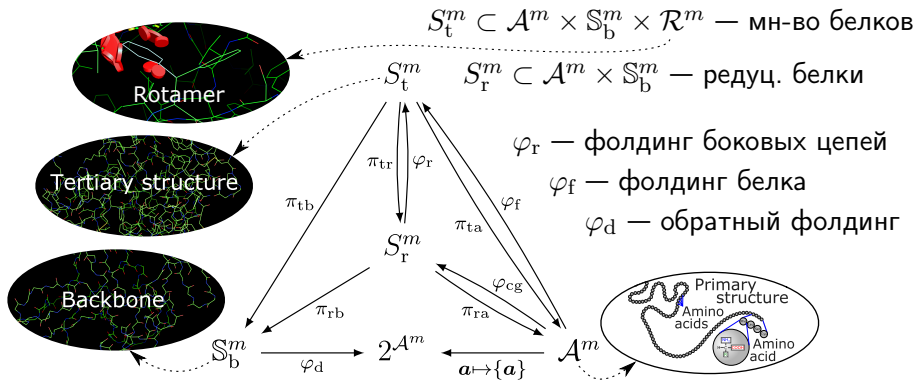
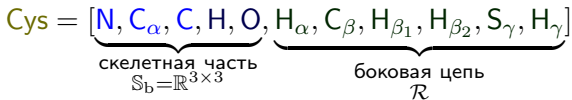
Vadim Strijov

Moscow Institute of Physics and Technology

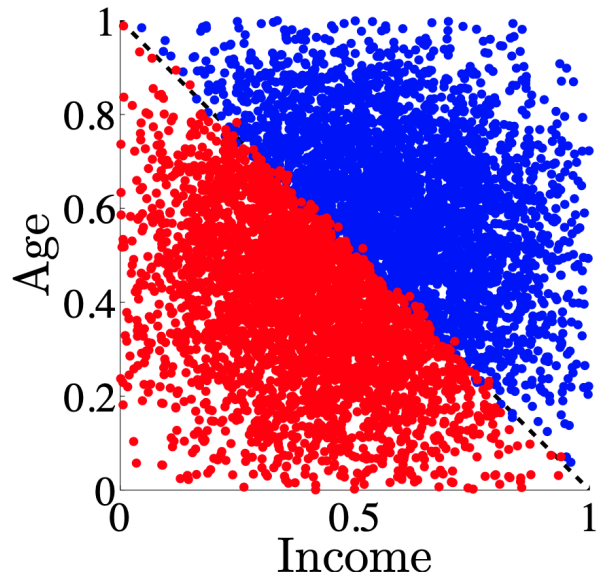
2021

Задачи структурной биологии для белков длины m

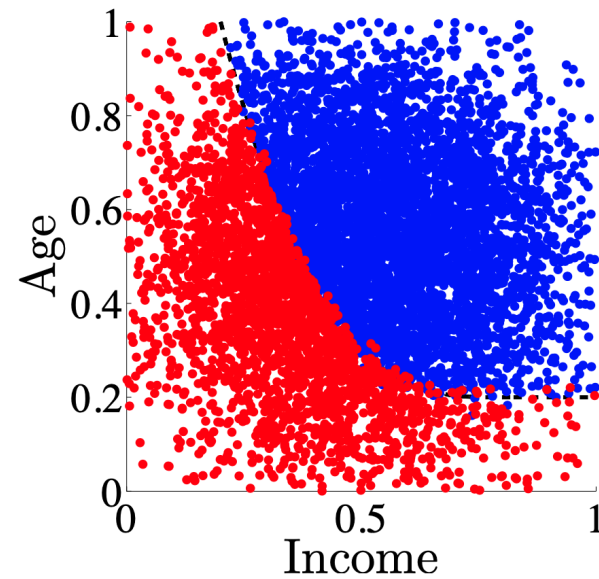
$$\mathcal{A} = \{\text{Ala, Arg, Asn, Asp, Cys, Glu, Gln, Gly, His, \dots, Trp, Tyr, Val}\}$$



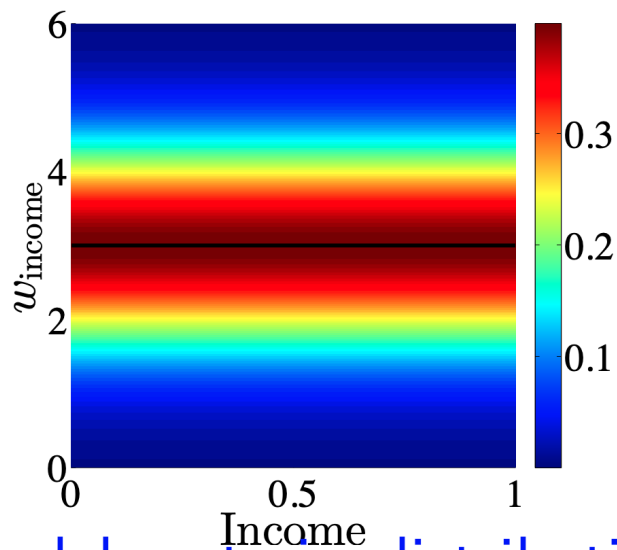
When one needs a multimodel?



Sample generation assumption

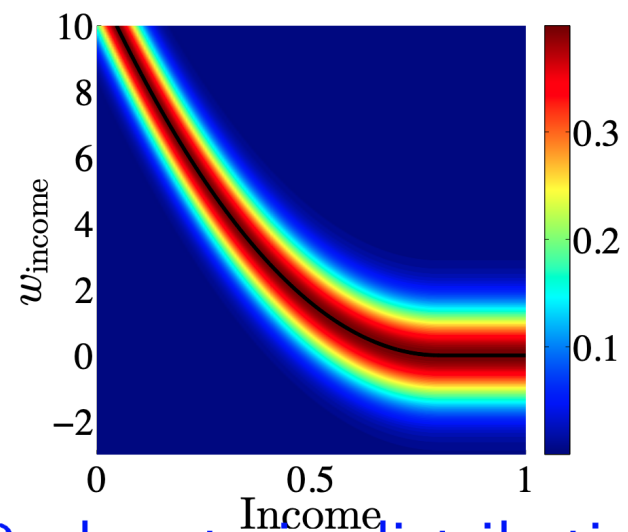


Real data



Model posterior distribution

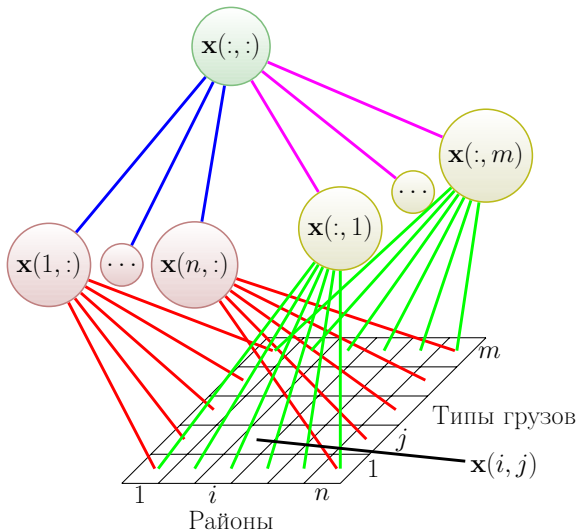
$$p(\mathbf{w}|\mathbf{X}, \mathbf{y})$$



Real posterior distribution

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y})$$

Условие согласованности прогнозов



$$x_t(:, :) = \sum_{i=1}^n x_t(i, :);$$

$$x_t(:, :) = \sum_{j=1}^m x_t(:, j);$$

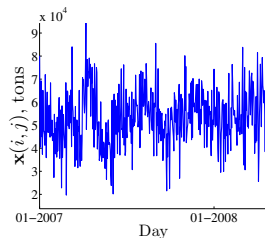
$$x_t(i, :) = \sum_{j=1}^m x_t(i, j),$$

$$i = 1, \dots, n;$$

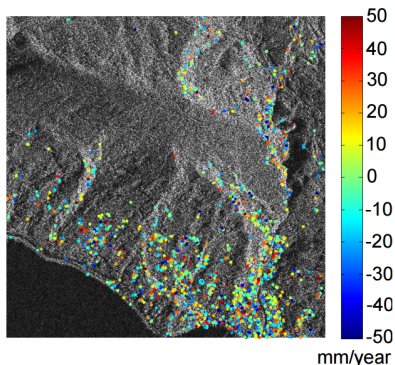
$$x_t(:, j) = \sum_{i=1}^n x_t(i, j),$$

$$j = 1, \dots, m;$$

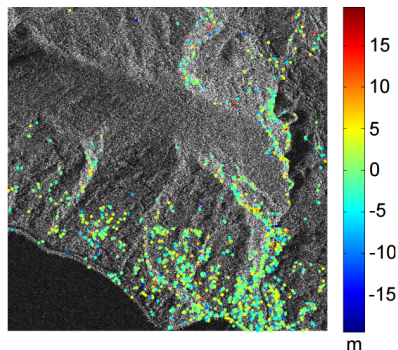
$$t = 1, \dots, T.$$



Анализ спутниковых снимков



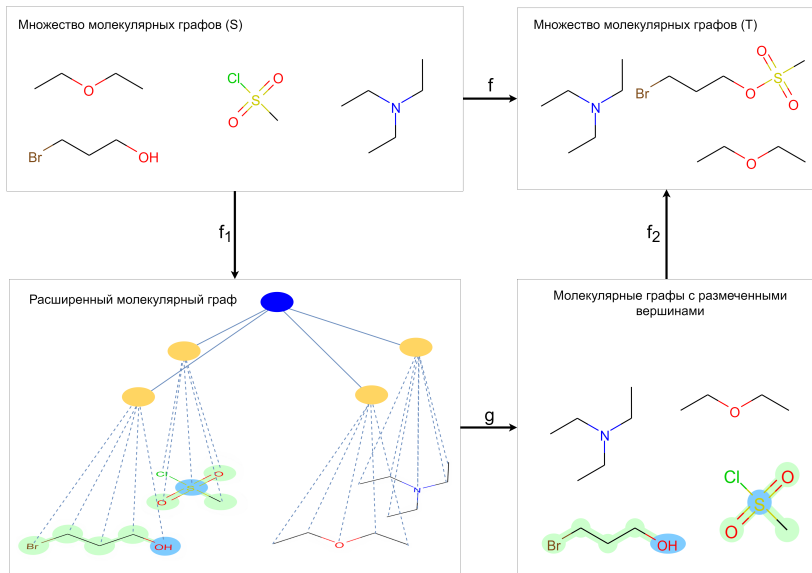
Скорость движения вдоль направления наблюдения



Погрешность в высоте

Рудаков К.В., Адуенко А.А., Рейер И.А., Василейский А.С., Карелов А.И., Стрижов В.В.
Алгоритмы выделения и смещения устойчивых отражателей на спутниковых снимках // Компьютерная оптика, 2015.

Структура решения



Задача регрессии в пространстве мультииндексных признаков

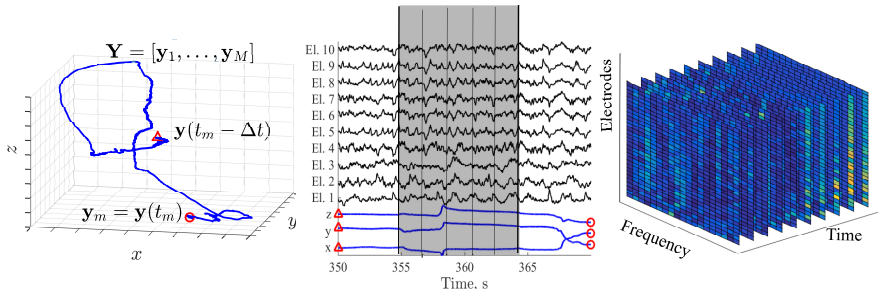
Решается задача восстановления регрессии

$$\hat{\mathbf{y}} = \sum_{i_1, \dots, i_D} [\underline{\mathbf{X}} * \underline{\mathbf{W}}]_{i_1 \dots i_D}, \quad \underline{\mathbf{W}} = \arg \min_{\mathbf{W}} \sum_{m=1}^M \|\underline{\mathbf{X}}_m * \underline{\mathbf{W}} - \mathbf{y}_m\|_F^2.$$

где $\mathbf{y} \in \mathbb{R}^n$, $\underline{\mathbf{X}}, \underline{\mathbf{W}} \in \mathbb{R}^{n_1 \times \dots \times n_D}$ — мультииндексные матрицы, $*$ обозначает операцию поэлементного умножения.

В векторном виде задача имеет вид $\mathbf{y} = \mathbf{x}^T \mathbf{w}$, где

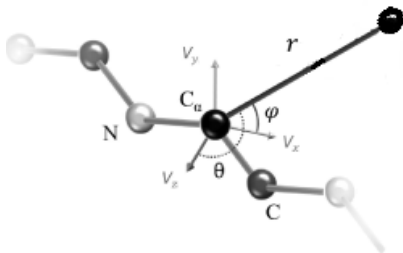
$$\mathbf{x} = \text{vec}(\underline{\mathbf{X}}) = x_{i_1(n_1-1) + \dots + i_{D-1}(n_{D-1}-1) + i_D} = [\underline{\mathbf{X}}]_{i_1 \dots i_D}.$$



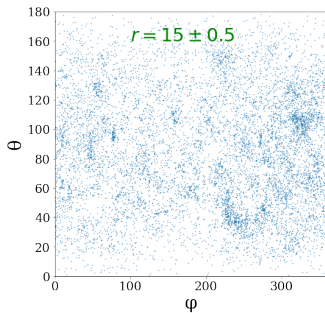
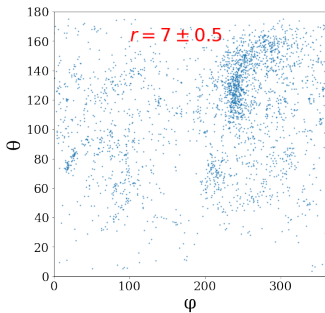
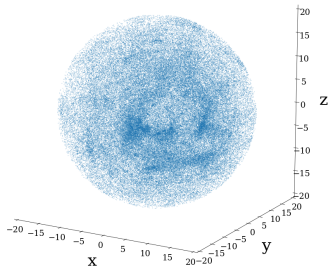
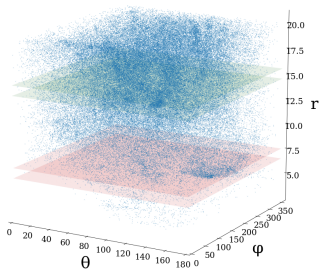
Описание молекулярной химической связи

В данной работе исследуются взаимные пространственные ориентации различных пар молекул, образующих между собой химическую связь. Эта связь характеризуется тремя параметрами:

- r — расстояние между молекулами, $r \in [3\text{\AA}, 20\text{\AA}]$;
- (θ, φ) — пара сферических углов, определяющих положение лиганда в системе координат аминокислоты, $\theta \in [0, \pi]$, $\varphi \in [0, 2\pi]$.



Представление выборки для пары ALA-C_{ar}



Сегментирование квазипериодического временного ряда

Квазипериодический временной ряд $\mathbf{s} = \{s(t_1), \dots, s(t_i), \dots, s(t_m)\}$ длины m определяется набором $\langle \mathbf{s}^*, a(i, s), f(i) \rangle$, так что

$$s(t_i) = a(i, s_{[f(i)]}),$$

где $\mathbf{s}^* = [s_1, \dots, s_T]^T$ — базовый сегмент,

$a(i, s), i \in \{1, \dots, m\}$ — трансформация формы

базового сегмента, $f(i) \mapsto \{1, \dots, T\}$ —

масштабирование по времени.

Теорема (Мотренко)

Для временного ряда \mathbf{s}

вида $s(t_i) = A_i \cos(2\pi w i + \phi)$ с $w \in (0, 1/2)$,

$\phi \in [0, 2\pi)$, $m \cdot w \in \mathbb{N}$ и $A_i : \exists C \in \mathbb{R} |A_i| < C \forall i$

главные компоненты \mathbf{y}_1 и \mathbf{y}_2 могут быть

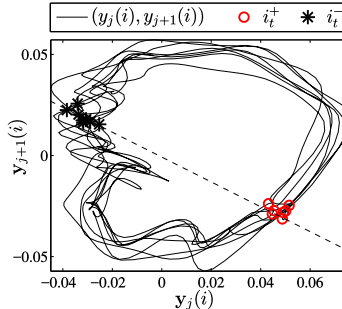
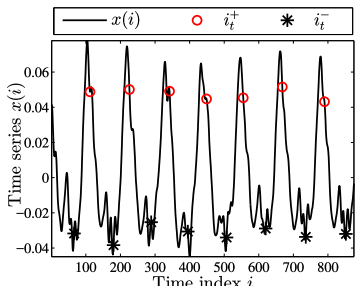
представлены в виде

$$y_1(l) = B_1(l) \cos(2\pi w l + \phi_1),$$

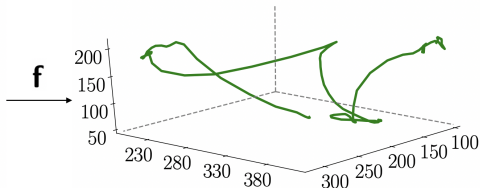
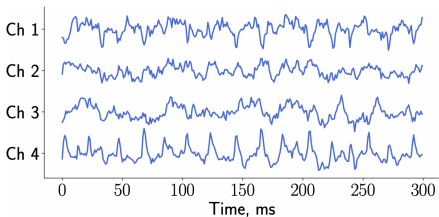
$$y_2(l) = B_2(l) \cos(2\pi w l + \phi_2),$$

$$\phi_1, \phi_2 \in [0, 2\pi), l = 1, \dots, m - N + 1$$

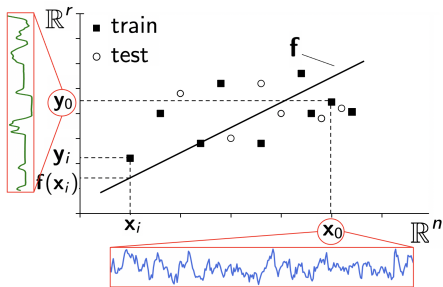
где разница между $|\phi_1 - \phi_2| \rightarrow \pi/2$.



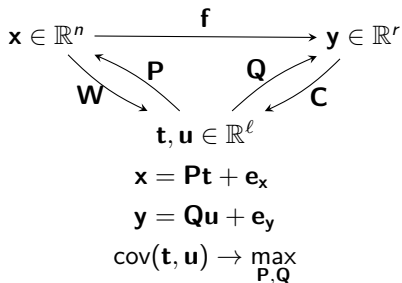
Восстановление зависимости в исходном и целевом пространствах



Прогностическая модель декодирования



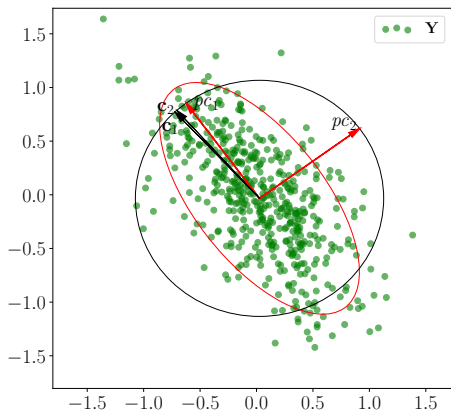
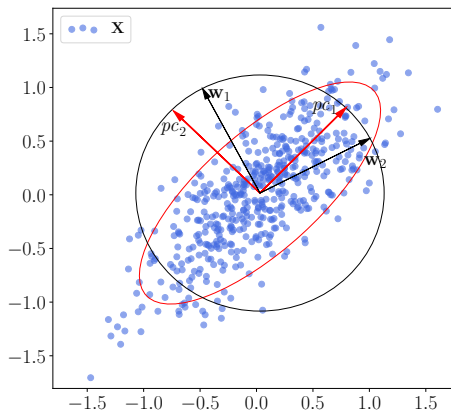
Согласование зависимостей в скрытом пространстве



Пример согласованной проекции в скрытое пространство

Исходные переменные $\mathbf{x}_i \sim \mathcal{N}(0, \Sigma)$.

Целевые переменные \mathbf{y}_i линейно зависят от pc_2 и не зависят от pc_1 .

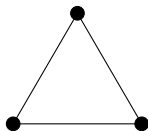


Согласование проекций матриц \mathbf{X} и \mathbf{Y} находит оптимальное скрытое представление, отклоняя вектора \mathbf{w}_k и \mathbf{c}_k от направления главных компонент.

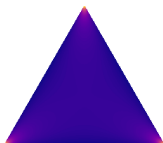
Априорное распределение на структуре модели

Каждая точка на симплексе задает модель.

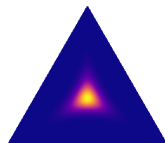
Распределение Гумбель-софтмакс: $\Gamma \sim \text{GS}(s, \lambda_{\text{temp}})$



$$\lambda_{\text{temp}} \rightarrow 0$$

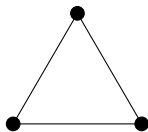


$$\lambda_{\text{temp}} = 0.995$$

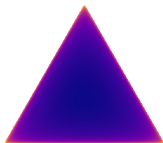


$$\lambda_{\text{temp}} = 5.0$$

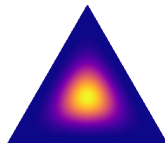
Распределение Дирихле: $\Gamma \sim \text{Dir}(s, \lambda_{\text{temp}})$



$$\lambda_{\text{temp}} \rightarrow 0$$



$$\lambda_{\text{temp}} = 0.995$$

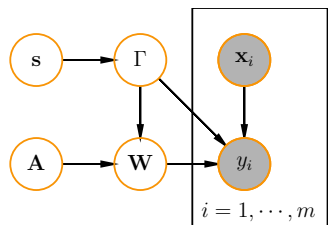


$$\lambda_{\text{temp}} = 5.0$$

Байесовский выбор модели

Базовая модель:

- параметры модели $\mathbf{w} \sim \mathcal{N}(0, \alpha^{-1})$,
- гиперпараметры модели $\mathbf{h} = [\alpha]$.



Предлагаемая модель:

- параметры модели $\mathbf{w}_r^{j,k} \sim \mathcal{N}(0, \gamma_r^{j,k} (\mathbf{A}_r^{j,k})^{-1})$, $\mathbf{A}_r^{j,k}$ — диагональная матрица параметров, соответствующих базовых функций $\mathbf{g}_r^{j,k}$, $(\mathbf{A}_r^{j,k})^{-1} \sim \text{inv-gamma}(\lambda_1, \lambda_2)$,
- структурные параметры модели $\Gamma = \{\gamma^{j,k}, (j, k) \in E\}$, $\gamma^{j,k} \sim \text{GS}(s^{j,k}, \lambda_{\text{temp}})$,
- гиперпараметры модели $\mathbf{h} = [\text{diag}(\mathbf{A}), s]$,
- метапараметры $\lambda_1, \lambda_2, \lambda_{\text{temp}}$.

Верхняя оценка правдоподобия модели

Оптимальные факторы $q^*(\theta)$, $q^*(\mathbf{m}, \mathbf{V}, \alpha)$ имеют вид

$$\ln q^* = E[\ln p(\mathbf{Z}, \theta, \mathbf{m}, \mathbf{V}, \alpha)] + \text{const}$$

и не вычисляются аналитически, так как правдоподобие L в модели $p(\mathbf{Z}, \theta, \mathbf{m}, \mathbf{V}, \alpha)$ содержит сумму экспонент $g(\mathbf{s}_n)$ в знаменателе,

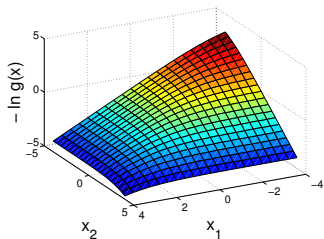
$$L(\mathbf{Z}|\theta, \alpha) = \prod_{n=1}^N \prod_{k=1}^{K_h} \left[\frac{\exp(s_{n,k})}{g(\mathbf{s}_n)} \right]^{z_{nk}}, \quad g(\mathbf{s}_n) = \sum_{k=1}^{K_h} \exp(s_{n,k}).$$

$\ln \frac{1}{g(\mathbf{x})}$ – вогнутая функция. Касательная плоскость к ней через точку ξ :

$$\zeta(\mathbf{x}, \xi) = -\ln(g(\xi)) - \nabla \ln(g(\xi))^T (\mathbf{x} - \xi).$$

Верхняя оценка правдоподобия L

$$\frac{1}{g(\xi)} \exp \left(s_{n,k} + \sum_{k'=1}^{K_h} \frac{\exp(\xi_{k'})}{g(\xi)} (\xi_{k'} - s_{n,k'}) \right).$$



Анализ качества предлагаемого метода выбора признаков

Решена задача прогнозирования многомерных временных рядов $\mathbf{y}(t) \in \mathbb{R}^3$ координат конечности по интервалам $\mathbf{s}(t - \Delta t)$ многомерных временных рядов $\mathbf{s}(t) \in \mathbb{R}^{N_{\text{ch}}}$ многоканальных электрокортикограмм.

Признаки:

$$\underline{\mathbf{X}}_m \in \mathbb{R}^{F \times N_{\text{ch}}}, \quad \underline{\mathbf{X}}_{mjn} = \begin{cases} s_n(t_m + \tau), & j = 1, \\ W_{mjn}, & j = 2, \dots, F + 1, \end{cases}$$

Прогноз в точке t_m :

$$\hat{\mathbf{y}}_m = \text{vec}(\underline{\mathbf{X}}_m)^T \hat{\mathbf{w}}.$$

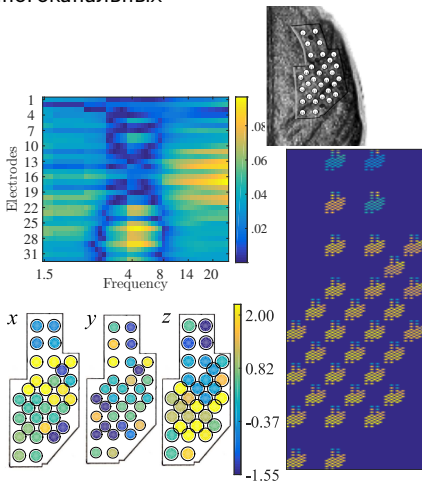
Качество прогнозирования:

коэффициент корреляции между $\hat{\mathbf{Y}}$ и \mathbf{Y} ,

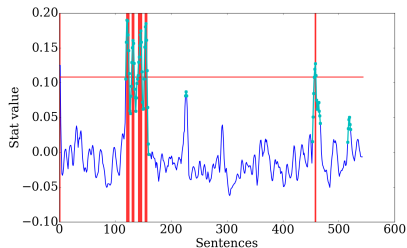
$$\text{corr}(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{\text{cov}(\hat{\mathbf{y}}, \mathbf{y})}{\sqrt{\text{cov}(\hat{\mathbf{y}}, \hat{\mathbf{y}})\text{cov}(\mathbf{y}, \mathbf{y})}}.$$

масштабированная ошибка MSE,

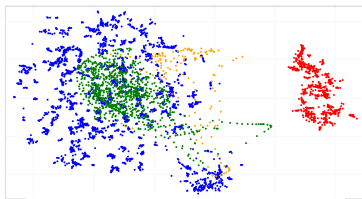
$$\text{sMSE}(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{\sum_{m=1}^M \|\hat{\mathbf{y}}_m - \mathbf{y}_m\|_2}{\sum_{m=1}^M \|\bar{\mathbf{y}} - \mathbf{y}_m\|_2}.$$



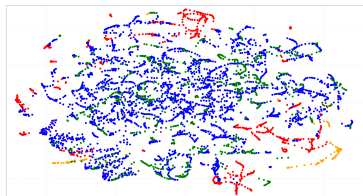
Снижение размерности с сохранением локальной структуры близости



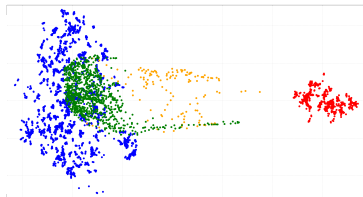
Построение признакового пространства для ряда s



Модифицированный t-SNE, $\mu = 10$



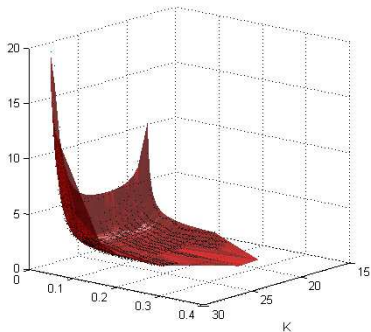
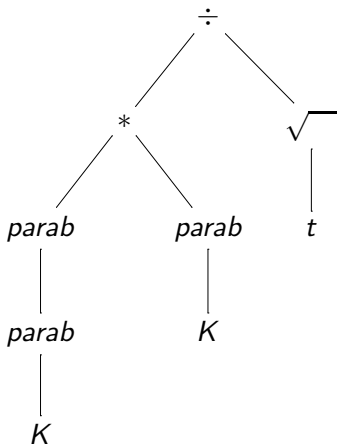
Исходный t-SNE



Модифицированный t-SNE, $\mu = 100$

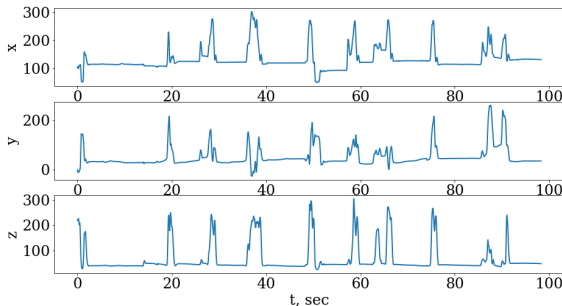
Результаты для опциона CLZ

$$\sigma = \frac{(w_1 K + w_2)(w_1 K^2 + w_2 K + w_3)^2}{\sqrt{t}}$$



Декодируемые сигналы электрокортикограммы

- Сигналы $\mathbf{s}(t) \in \mathbb{R}^{N_{ch}}$. N_{ch} – число электродов
- Координаты электродов $\mathbf{Z} = \{(\mathbf{z}_j \in \mathbb{R}^2, j \in \{1 \dots, N_{ch}\})\}$
- Положение кисти в пространстве $\mathbf{y}(t) \in \mathbb{R}^3$



Координата руки



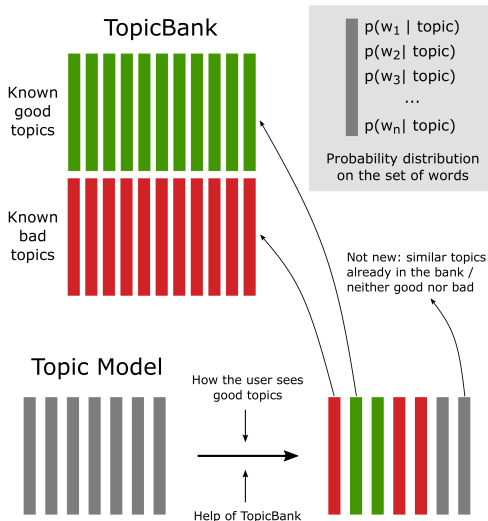
Пространственное
расположение
электродов

Chao ZC, Nagasaka Y, Fujii N (2010). "Long-term asynchronous decoding of arm motion using electrocorticographic signals in monkeys." *Frontiers in Neuroengineering* 3:3.

Банк тем: сохранение интерпретируемых тем

Банк тем — модель
полного набора тем:
таких тем, которые

- 1) интерпретируемы,
- 2) существенно
различны,
- 3) обеспечивают
высокое
правдоподобие
модели
 $p(\Phi, \Theta | D)$.

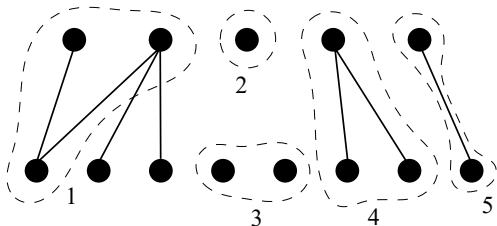


Построение банка тем

Аналогично построению двухуровневой иерархической тематической модели:

$$\underbrace{p(w | t)}_{\varphi_{wt}^{parent}} = \sum_{s \in S} \underbrace{p(w | s)}_{\varphi_{ws}^{child}} \underbrace{p(s | t)}_{\psi_{st}} \quad \text{Hierarchy}$$

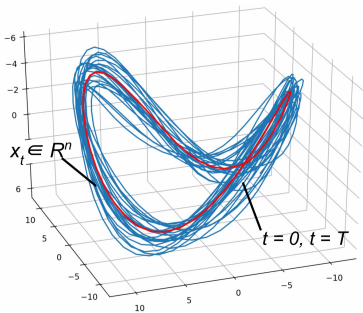
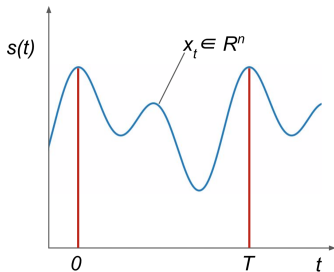
$$\underbrace{p(w | t)}_{\varphi_{wt}^{bank}} = \sum_{s \in S} \underbrace{p(w | s)}_{\varphi_{ws}^{new}} \underbrace{p(s | t)}_{\psi_{st}} \quad \text{TopicBank}$$



№	Hierarchy	TopicBank
1	ok	no
2	ok	ok
3	no	ok
4	ok	maybe
5	ok	maybe

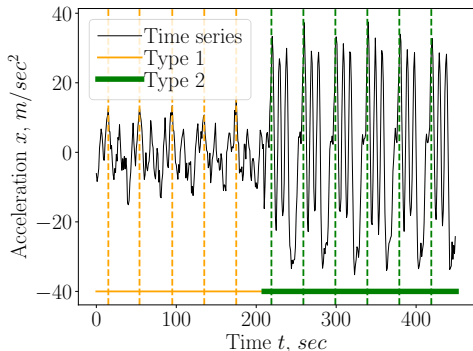
Фазовая траектория

На рисунке представлен временной ряд и проекция его фазовой траектории в трехмерное пространство. $\mathbf{x}_t = \mathbf{x}(t)$ – точка на фазовой траектории в момент времени t .

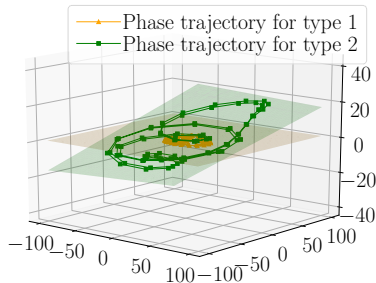


Сегмент — последовательность точек временного ряда, которая относится к одному характерному физическому действию человека: шаг, прыжок.

Цепь — последовательность сегментов, которые образуют квазипериодическую последовательность точек.



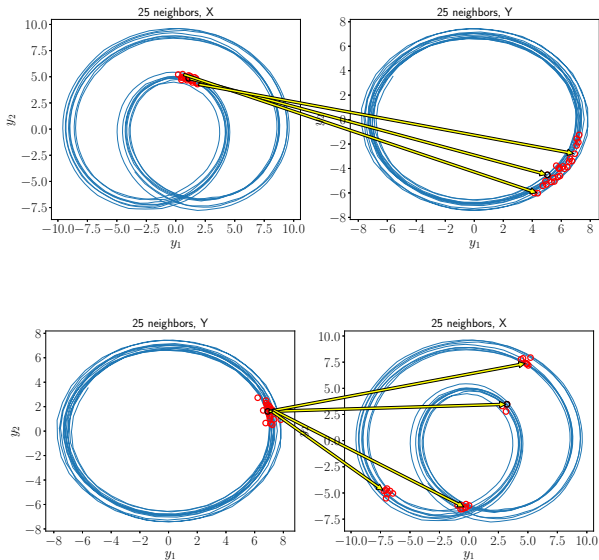
(a)



(b)

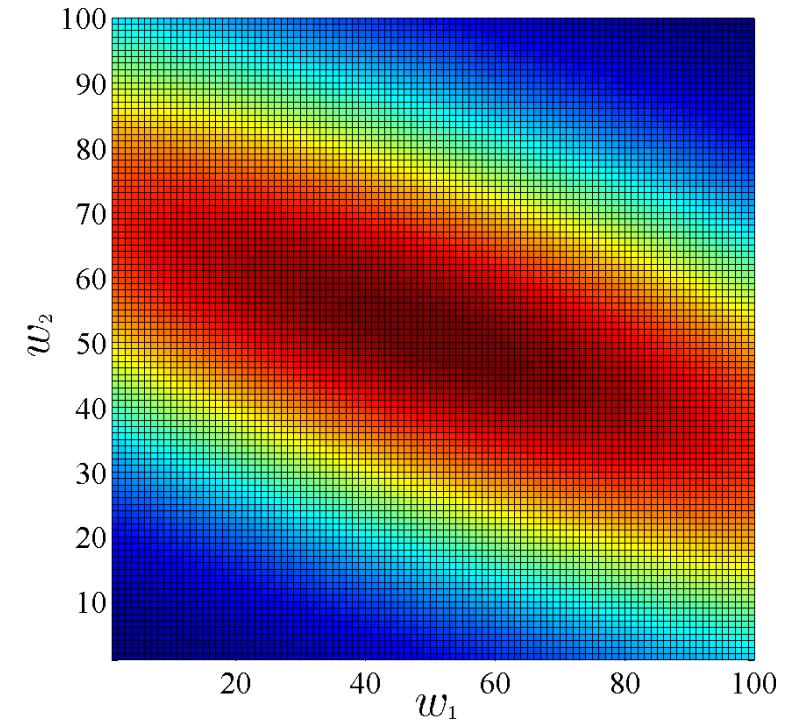
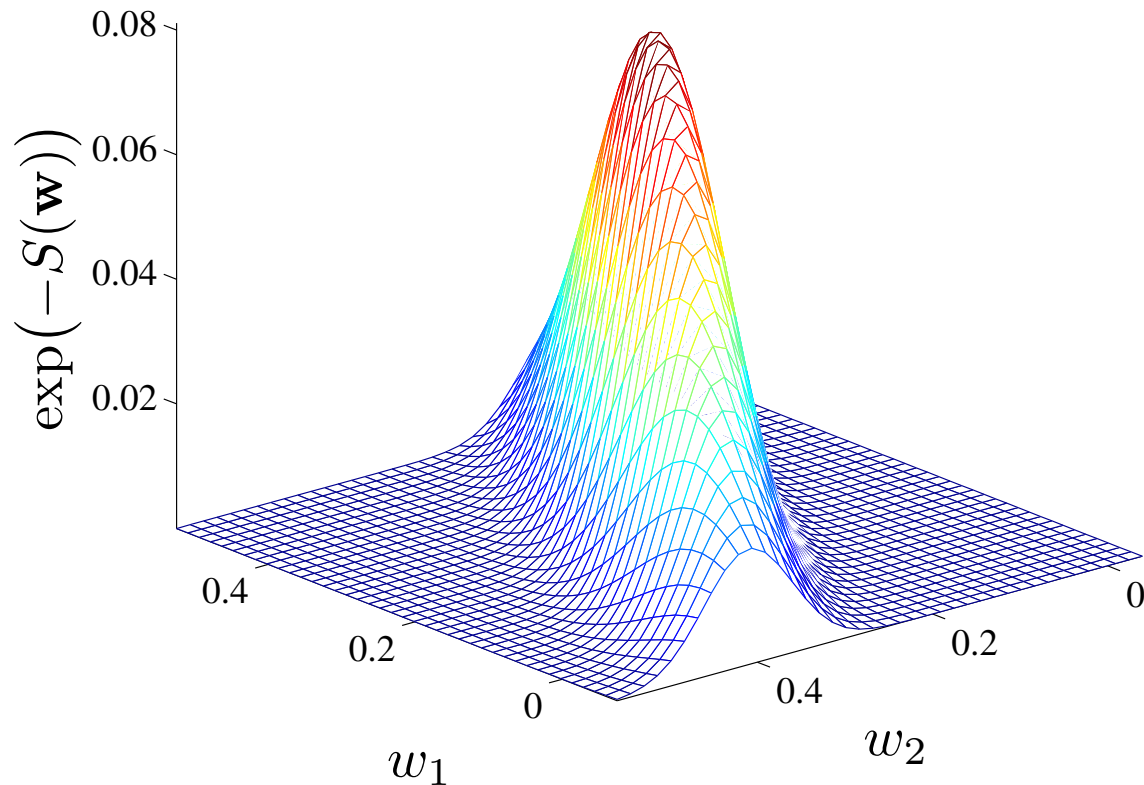
а) временной ряда разбитый на сегменты; б) проекции на плоскость фазовых траекторий временного ряда, которые относятся к Type 1 и Type 2.

Ближайшие соседи на фазовых траекториях



Empirical distribution of model parameters

The value of error function $S(\mathbf{w}|\mathcal{D}, f)$ depends on parameters.



x-axis and y-axis: parameters \mathbf{w} , z-axis: $\exp(-S(\mathbf{w}))$

Probabilistic model selection

Bayesian inference delivers the error function $S(\mathbf{w})$

$$p(\mathbf{w}|\mathcal{D}, \mathbf{A}, \mathbf{B}, \mathbf{f}) = \frac{p(\mathcal{D}|\mathbf{w}, \mathbf{B}, \mathbf{f})p(\mathbf{w}|\mathbf{A}, \mathbf{f})}{p(\mathcal{D}|\mathbf{A}, \mathbf{B}, \mathbf{f})}.$$

Posterior Likelihood Prior
Evidence
(to select a model)

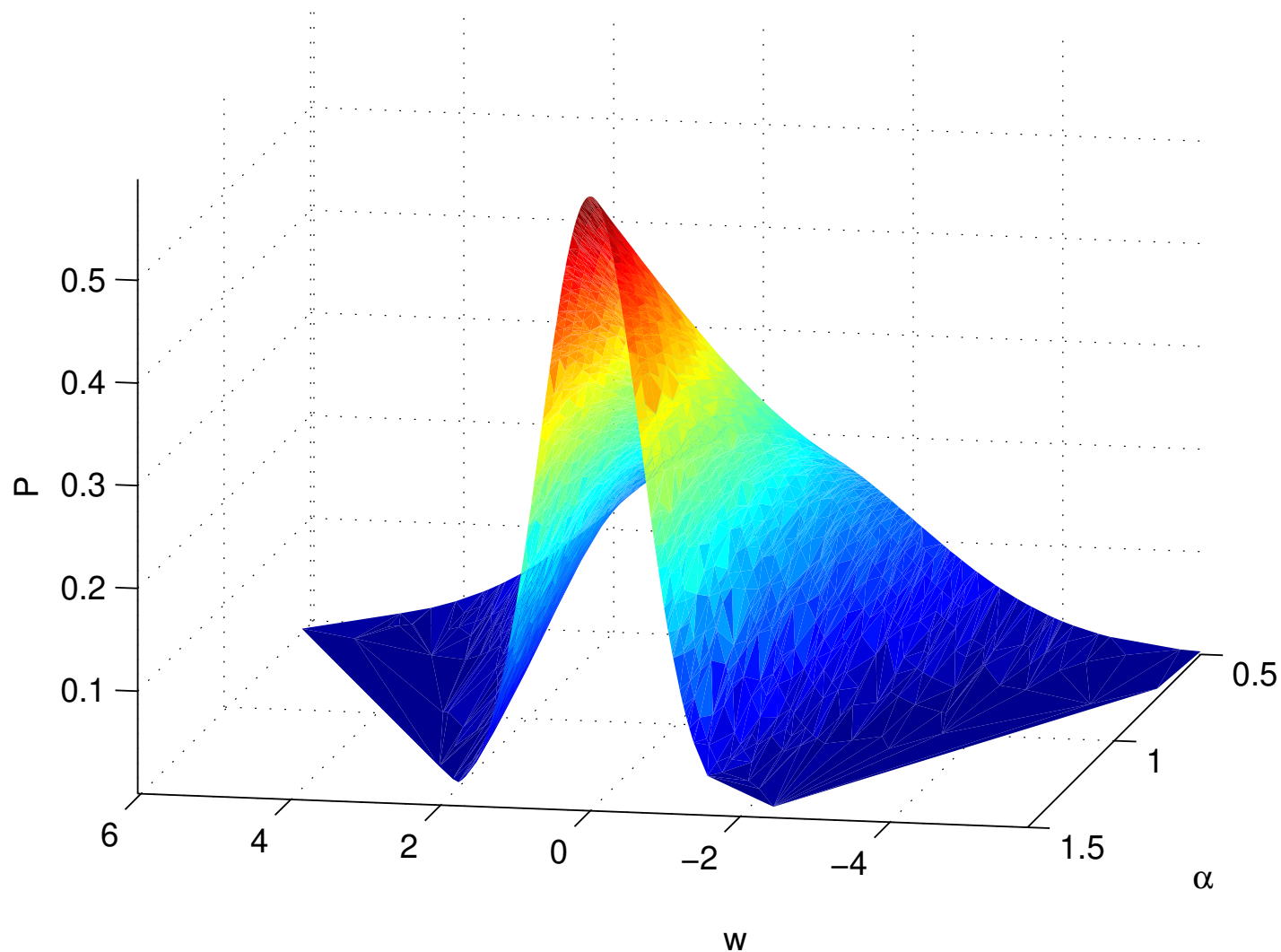
Write the error function given hyperparameters \mathbf{A}, \mathbf{B}

$$S(\mathbf{w}) = \underbrace{\frac{1}{2}(\mathbf{y} - \mathbf{f})^T \mathbf{B}(\mathbf{y} - \mathbf{f})}_{\text{approximation error}} + \underbrace{\frac{1}{2}(\mathbf{w} - \hat{\mathbf{w}})^T \mathbf{A}(\mathbf{w} - \hat{\mathbf{w}})}_{\text{regularisation error}},$$

$$S = E_D + E_w = \lambda^T \mathbf{s}, \quad \text{metaparameters } \lambda = \frac{1}{2}.$$

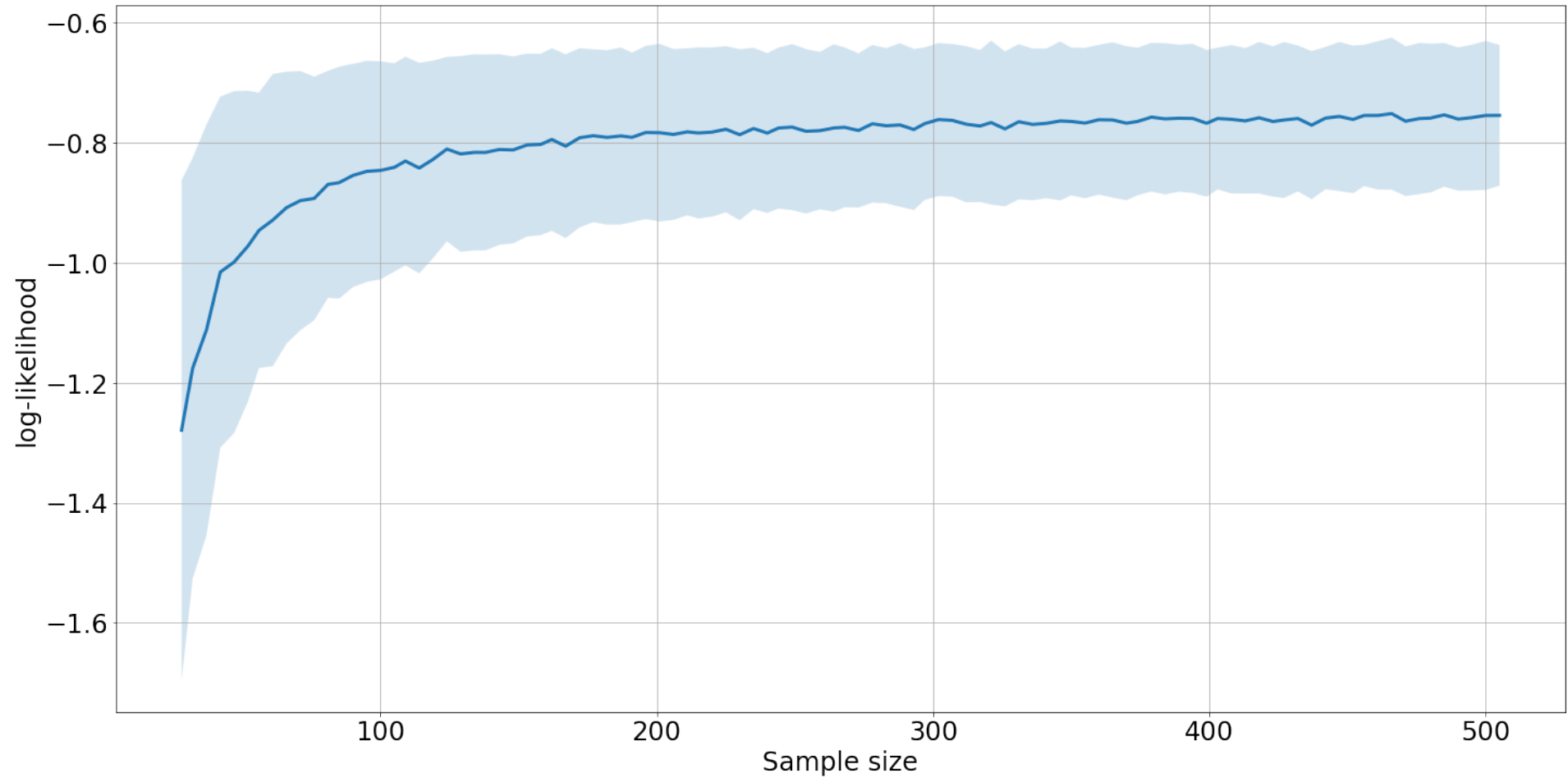
Evidence of the model

depends on both, error E_D (likelihood) and regularisation E_w (prior).

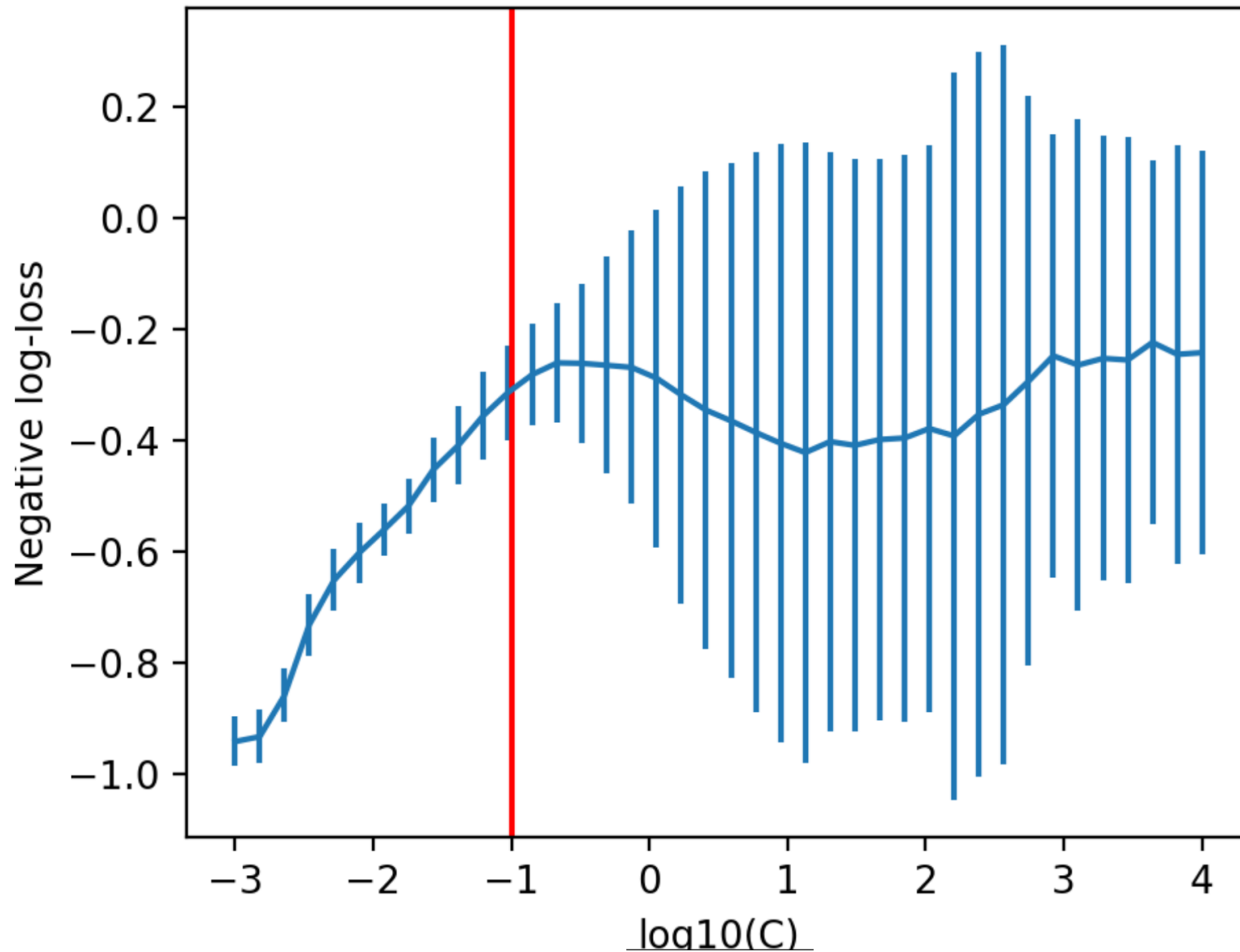


Parameters w , variance α^{-2} , and $p(w|\mathcal{D}, \alpha)$ is the evidence.

– Error and its variance for a reinforced sample set



Variance of error increasing over model complexity

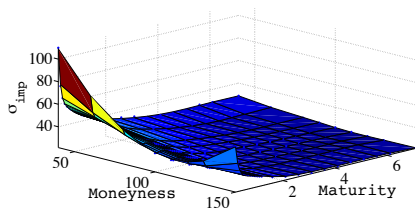
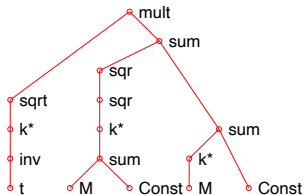


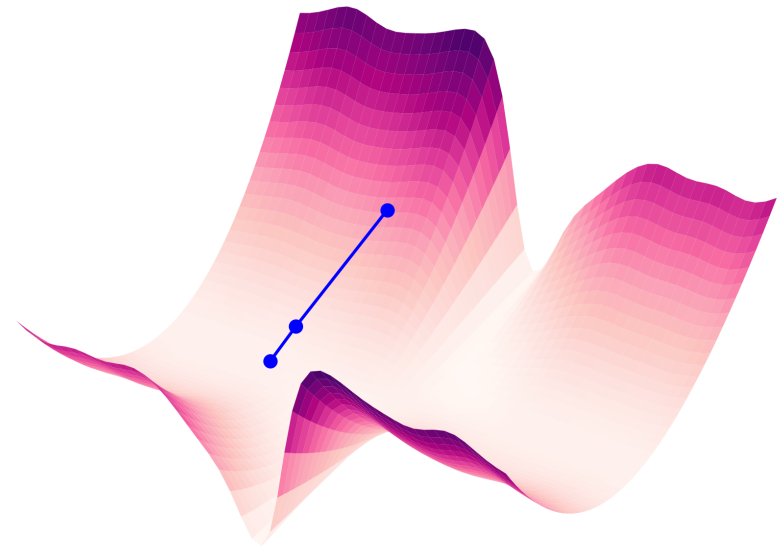
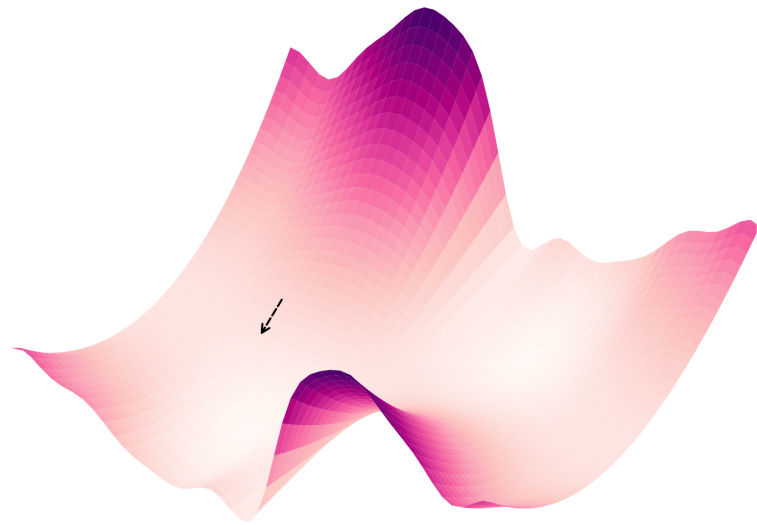
Complexity, number of model parameters

Resulting models

Resulting model

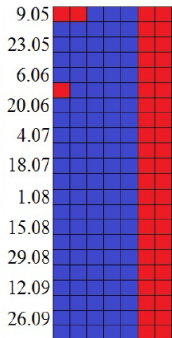
$$\sigma_{imp} = \frac{(k_2 M + k_3)^4 + k_4 M + c}{\sqrt{k_1 t}}$$



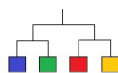
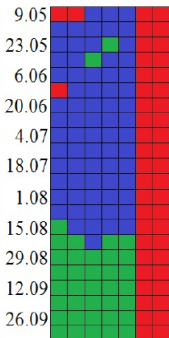




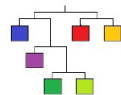
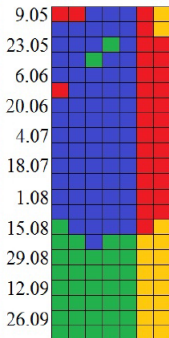
Пн Bc



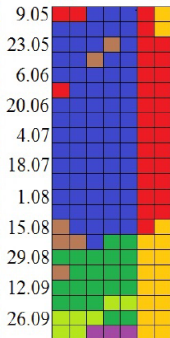
Пн Bc

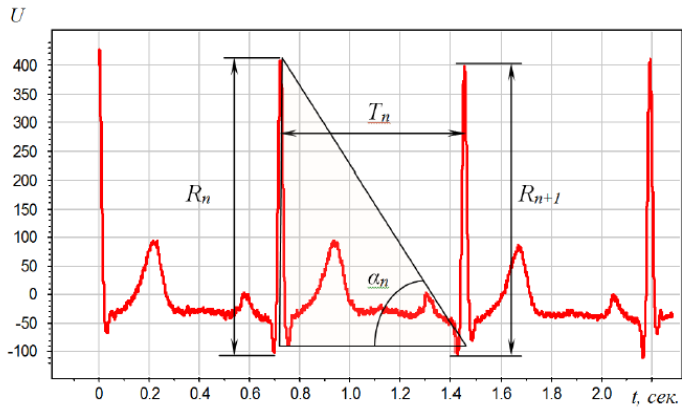


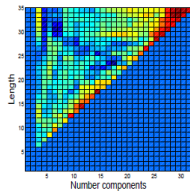
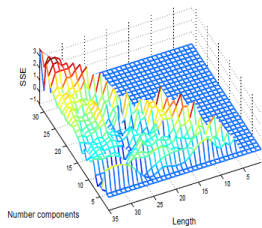
Пн Bc

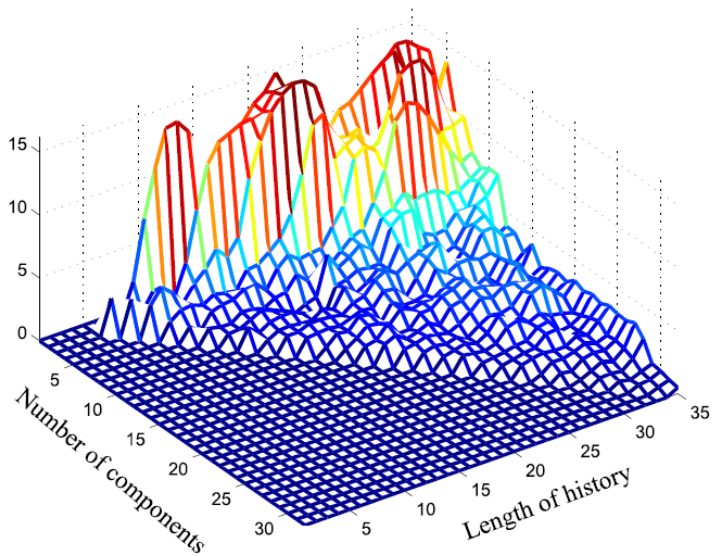


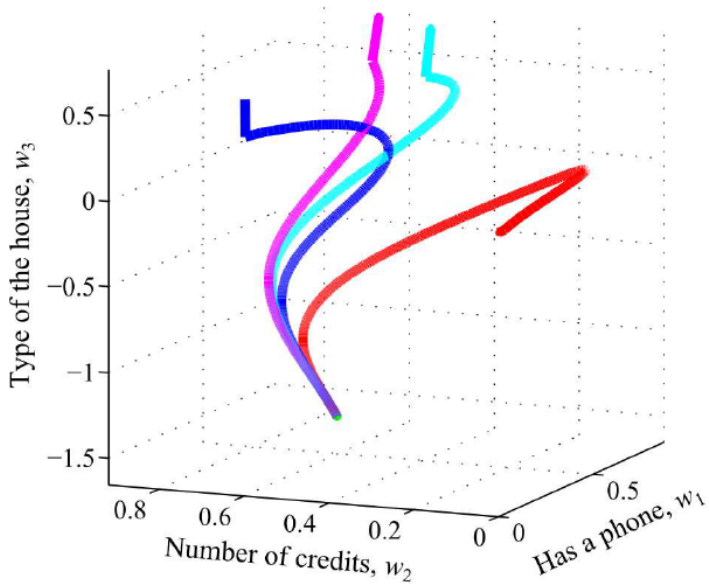
Пн Bc

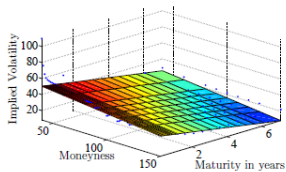
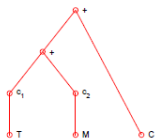




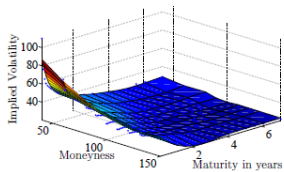
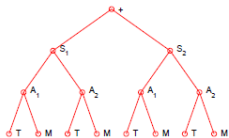




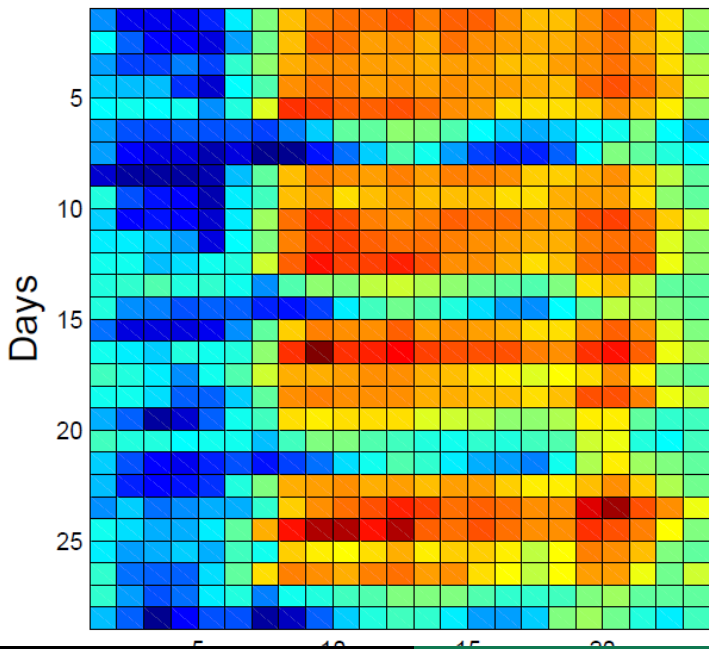


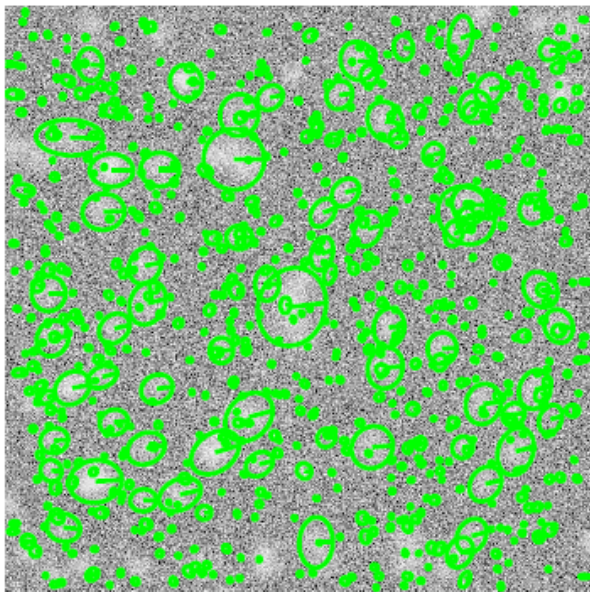


a)

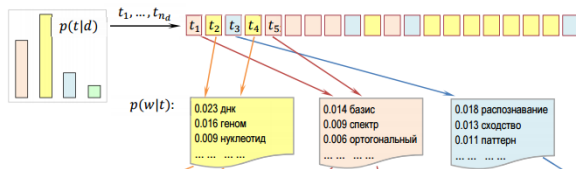


b)



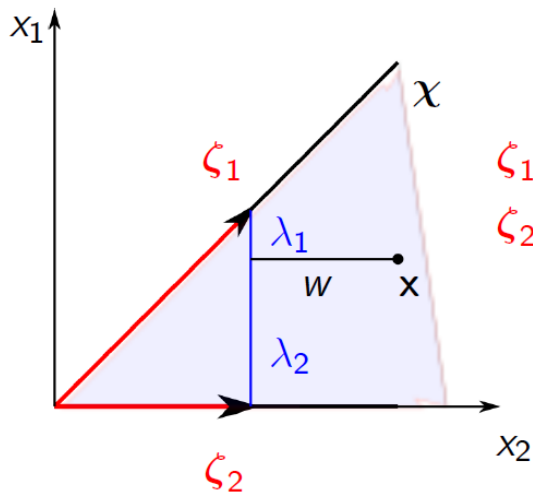


Тематическое моделирование и матричное разложение



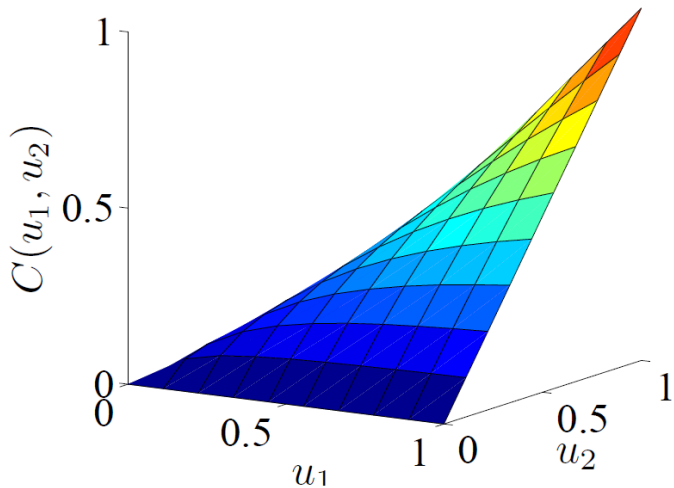
w_1, \dots, w_{n_d} :

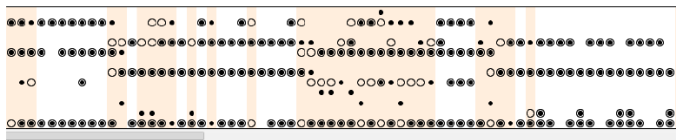
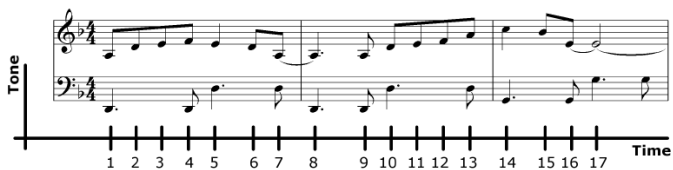
Разработан спектрально-аналитический подход к выявлению размытых протяженных повторов в геномных последовательностях. Метод основан на разномасштабном оценивании сходства нуклеотидных последовательностей в пространстве коэффициентов разложения фрагментов кривых GC- и GA-содержания по классическим ортогональным базисам. Найдены условия оптимальной аппроксимации, обеспечивающие автоматическое распознавание повторов различных видов (прямых и инвертированных, а также тандемных) на спектральной матрице сходства. Метод одинаково хорошо работает на разных масштабах данных. Он позволяет выявлять следы сегментных дупликаций и мегасателлитные участки в геноме, районы синтении при сравнении пары геномов. Его можно использовать для детального изучения фрагментов хромосом (поиска размытых участков с умеренной длиной повторяющегося паттерна).



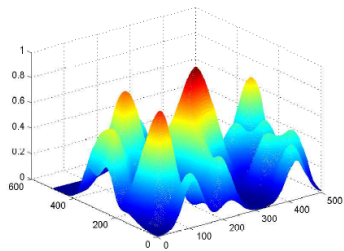
$$\zeta_1 = [1; 1],$$

$$\zeta_2 = [1; 0]$$

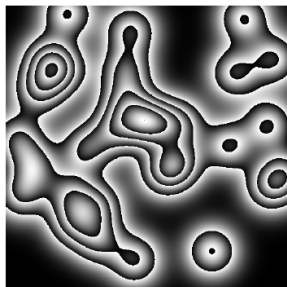




Пустые круги — истинные полутона, точки — предсказанные, ошибки подсвечены, горизонтальная ось — время.



(a) Высота



(б) Фазовая составляющая без учета шума

